

**Faculdade de Engenharia da Universidade do Porto**



## **Análise e Classificação de Imagem Hiper-espectral**

**Borgine Vasques Gurué**

Dissertação realizada no âmbito do  
Mestrado em Engenharia da Informação

Orientador: Prof. Dr. Jaime dos Santos Cardoso

Julho de 2017



© Borgine Vasques Gurué, 2017

# **Análise e Classificação de Imagem Hiper-espectral**

**Borgine Vasques Gurué**

Mestrado em Engenharia da Informação

Aprovado em provas públicas pelo Júri:

Presidente: Prof. Dr. António Pedro Rodrigues Aguiar

Vogal Externo: Doutora Inês Campos Monteiro Sabino Domingues

Orientador: Prof. Dr. Jaime dos Santos Cardoso

---

Julho de 2017

# Resumo

A incidência de incêndios florestais tem sido um fenómeno devastador para várias comunidades e contribui em alta para a destruição do próprio ecossistema e degradação das áreas vegetacionais. Este trabalho tem como propósito desenvolver modelos para a identificação de áreas propensas a ocorrência de incêndios e que, de certa forma facilite o mapeamento dessas zonas, de modo a manter as mesmas em alerta máximo nos períodos em que se verificam os maiores índices de queimadas. Os modelos devem distinguir e prever áreas com diferentes riscos de ocorrência de incêndios, entre eles, risco baixo, médio e alto. Com estas informações, estes modelos facilitam a tomada de decisão, a prevenir futuras ocorrências de incêndios nas zonas identificadas como de alto risco, durante um período específico do ano (os índices mais altos ocorrem sempre em períodos secos).

O estudo baseou-se em aplicar técnicas e modelos/algoritmos de *machine learning*, tanto para análise dos dados como para classificação, de acordo com as classes predefinidas, sobre cobertura vegetal no distrito de Mandimba, província do Niassa - Moçambique.

As experiências foram baseadas em quatro algoritmos diferentes, dos quais, os primeiros três considerados tradicionais (*Random Forest*, *Gradient Boosting*, *Support Vector Machine*) e um avançado (*Deep Belief Networks* - módulo de *Deep Learning*).

Os resultados mostram que, embora os algoritmos tradicionais consigam dar bons resultados, o *deep learning*, quando bem implementado/aplicado, dá-nos resultados ainda melhores.

Palavras-chave: Análise de imagens hiper-espectrais, *machine learning*, classificação.



# Abstract

The incidence of forest fires has been a devastating phenomenon for several communities and contributes to a sizable destruction of the ecosystem itself and the degradation of areas of fauna. This work aims to develop models to the identification of areas that are prone to the occurrence of fires and, makes it easier to map these areas, in order to keep them under maximum alert for periods for which there are higher indices of fires. The models should differentiate and predict areas of differing risk for the occurrence of fires, namely in low, medium and high risk. This information helps at decision making, predicting future occurrences of fires in areas identified as high risk, during a specific period of the year (higher indices are always related to the dry seasons).

The study was based in the application of techniques and models/algorithms of *machine learning*, both in terms of the data analysis for classification, according to predefined classes, over the fauna cover on the district of Mandimba, a province of Niassa - Mozambique.

Experiments were based on four different algorithms, of which, the first three were considered traditional ones (*Random Forest*, *Gradient Boosting*, *Support Vector Machine*) and an advanced one (*Deep Belief Networks* - a part of *Deep Learning*).

Results show that, although traditional algorithms are able to obtain good results, *deep learning*, when well implemented/applied, is able to produce even better results.

Keywords: Hyperspectral image processing, *machine learning*, classification.



# Agradecimentos

A Deus, pelo dom da vida.

Ao Camões IP - Instituto de Cooperação e Língua, por ter financiado os meus estudos.

Ao Professor Jaime dos Santos Cardoso, por ter aceite o meu desafio e orientado com muita paciência o trabalho e pelos ensinamentos por ele transmitidos que os levarei para o meu ambiente profissional.

Ao Professor António Pedro Rodrigues Aguiar, por ter aceite a minha candidatura, na qualidade de diretor do curso.

A minha família, pela força e pelo acompanhamento ao longo dessa jornada.

Ao Tiago Manuel Andrade dos Santos e sua equipa de pesquisa, pela disponibilização dos dados usados para materialização desta dissertação.

E por fim aos meus amigos de “trincheiras”, Rogers e Francelino pela força que sempre deram, aos meus “mestres” Ricardo Cruz e Kelvin Fernandes pelas instruções que me transmitiram, ao Cá Baltazar pelas discussões académicas que sempre tivemos e aos amigos “brazucas” Fernando e Mailson e suas parceiras, pelo acolhimento que deram numa cidade onde todos éramos estrangeiros, vai um muito obrigado a todos.





# Índice

<b>Capítulo 1 .....</b>	<b>15</b>
Introdução .....	15
1.1 - Enquadramento .....	15
1.2 - O Projeto .....	16
1.3 - Objetivos .....	16
1.4 - Área de estudo .....	17
1.5 - Contribuições .....	18
1.6 - Estrutura da dissertação .....	19
<b>Capítulo 2 .....</b>	<b>20</b>
Conceitos fundamentais e Revisão bibliográfica .....	20
2.1- Imagem hiper-espectral vs multiespectral .....	20
2.2- Análise de Imagem .....	21
2.2.1 Análise de Imagem Digital .....	21
2.3- Processamento Digital de Imagem .....	21
2.3.1 Aquisição da Imagem .....	22
2.3.2 Pré-processamento .....	23
2.3.3 Segmentação .....	23
2.3.3.1 Morfologia Matemática .....	25
2.3.3.1.1 Conceitos básicos sobre teoria de conjuntos .....	25
2.3.4 Extração de Atributos .....	27
2.3.5 Classificação e Reconhecimento .....	28
2.3.5.1 Random Forest Classifier .....	29
2.3.5.2 Gradient Boosting Classifier .....	29
2.3.5.2.1 Algoritmo gradient boosting .....	30
2.3.5.3 Deep Learning Classifier .....	30
2.3.5.4 Máquina de Vector de Suporte .....	32
2.4- Revisão Bibliográfica .....	33
<b>Capítulo 3 .....</b>	<b>37</b>
Metodologia & Resultados .....	37
3.1- Introdução .....	37
3.2- Ferramentas utilizadas .....	38
3.3- Aquisição de dados .....	38
3.4- Experiência e resultados com modelos propostos .....	42
3.5- Melhoria dos resultados .....	43
3.5.1 Agregação de classes .....	44
3.5.2 Combinação de classificadores (ensemble) .....	45
3.6- Uso de dados com diferença temporal (2002 vs 2005) .....	45
3.6.1 Estratégia ensemble .....	48
3.6.2 Comparação dos algoritmos GB, RF e SVM .....	49

3.6.3 Visualização da previsão das classes na imagem .....	50
3.7- Entendendo o Amazonas pelo espaço .....	53
<b>Capítulo 4 .....</b>	<b>56</b>
Conclusão .....	56
<b>Referências .....</b>	<b>58</b>

# Lista de figuras

FIGURA 1-1 MAPA DE MOÇAMBIQUE A), PROVÍNCIA DO NIASA B) COM DESTAQUE O DISTRITO DE MANDIMBA .....	18
FIGURA 2-1 PRINCIPAIS ETAPAS DE UM SISTEMA DE PROCESSAMENTO DIGITAL DE IMAGEM, ADAPTADO DE [12] .....	22
FIGURA 2-2 APLICAÇÃO DO FILTRO DE MEDIANA PARA REMOÇÃO DE RUÍDO [15] .....	23
FIGURA 2-3 (A) IMAGEM ORIGINAL, (B) SEGMENTAÇÃO POR BINARIZAÇÃO E (C) SEGMENTAÇÃO POR DETECÇÃO DE BORDAS [17] .....	24
FIGURA 2-4 A) REPRESENTA DOIS CONJUNTOS A E B. B) UNIÃO ENTRE A E B, DENOTA-SE POR $A \cup B$ . C) INTERSECÇÃO DE A E B, DENOTA-SE POR $A \cap B$ . D) DIFERENÇA DE A E B, DENOTA-SE POR $A - B$ . E) COMPLEMENTO DE A, DENOTA-SE POR (A)C, FIGURA ADAPTADA DE [11]. .....	25
FIGURA 2-5 A) TRANSLAÇÃO DE A POR Z. B) REFLEXÃO DE B .....	26
FIGURA 2-6 REPRESENTAÇÃO GERAL PARA UM PROCESSO DE CLASSIFICAÇÃO, ADAPTADA DE [13]. .....	28
FIGURA 2-7 - ARQUITETURA DEEP LEARNING [33] .....	31
FIGURA 2-8 - EXEMPLO DE SVM (SEPARAÇÃO POR HÍPER-PLANOS) [38] .....	32
FIGURA 3-1 VISUALIZAÇÃO DE OBSERVAÇÕES POR CLASSES .....	46
FIGURA 3-2 ARQUITETURA CONCEPT DRIFT, ADAPTADA DE [53] .....	48
FIGURA 3-3 COMPARAÇÃO POR ACURÁCIA MÉDIA DOS ALGORITMOS GB, RF E SVM .....	50
FIGURA 3-4 IMAGEM ORIGINAL REPRESENTANDO A ÁREA EM ESTUDO, SEM .....	51
FIGURA 3-5 PREVISÕES DAS CLASSES COM RANDOM FOREST .....	51
FIGURA 3-6 PREVISÕES DAS CLASSES COM GRADIENT BOOSTING .....	52
FIGURA 3-7 PREVISÕES DAS CLASSES COM SVM .....	52
FIGURA 3-8 IMAGEM DE TESTE (TEST_7) SOBRE UMA PARTE DO AMAZONAS. PREVISÕES: PRIMARY, ROAD, CLEAR, AGRICULTURE, HABITATION .....	55

## Lista de tabelas

TABELA 2-1 OS TRÊS OPERADORES LÓGICOS BÁSICOS.....	26
TABELA 2-2 - RESUMO DE RESULTADOS SOBRE ESTUDOS PASSADOS.....	36
TABELA 3-1 – DESCRIÇÃO DAS BANDAS DO SATÉLITE LANDSAT 7 .....	39
TABELA 3-2 - IMAGENS DE CADA BANDA ESPECTRAL: A) BANDA 2, B) BANDA 3, C) BANDA 4, D) BANDA 5), E) BANDA 7, F) NDVI E G) VI7.....	40
TABELA 3-3 - DESCRIÇÃO DAS CLASSES PROPOSTAS .....	41
TABELA 3-4 - NÚMERO DE OBSERVAÇÕES POR CLASSE (DADOS 2002) .....	42
TABELA 3-5 – PRECISÃO/ACCURACY POR CLASSIFICADOR.....	42
TABELA 3-6 – PRECISÃO/ACCURACY DO CLASSIFICADOR POR CADA CLASSE.....	43
TABELA 3-7 MELHORIA DE RF E GB POR AJUSTE DE PARÂMETROS.....	43
TABELA 3-8 AGREGAÇÃO DE CLASSES EM TRÊS NÍVEIS .....	44
TABELA 3-9 ACURÁCIA COM CLASSES AGREGADAS.....	44
TABELA 3-10 COMBINAÇÃO DOS CLASSIFICADORES RF, GB E SVM.....	45
TABELA 3-11 CLASSIFICAÇÃO COM DADOS COM DIFERENÇA TEMPORAL (2002 VS 2005) .....	46
TABELA 3-12 CLASSIFICAÇÃO COM DADOS COM DIFERENÇA TEMPORAL (2005 VS 2002) .....	47
TABELA 3-13 RESULTADO DO ENSEMBLE COM DADOS DIFERENTES NO TEMPO .....	49
TABELA 3-14 COMPARAÇÃO POR ACURÁCIA MÉDIA DOS ALGORITMOS GB, RF E SVM.....	49

# Abreviaturas e Símbolos

## Lista de abreviaturas

DBN	<i>Deep Belief Network</i>
DL	<i>Deep Learning</i>
GB	<i>Gradient Boosting</i>
NDVI	<i>Normalized Difference Vegetation Index</i>
OBIA	<i>Object-based Image Analysis</i>
RF	<i>Random Forest</i>
SVM	<i>Support Vector Machine</i>
TIFF	<i>Tagged Image File Format</i>

# Capítulo 1

## Introdução

### 1.1 - Enquadramento

O uso de imagens na identificação de padrões semelhantes numa dada característica específica, tem sido objeto de estudo em áreas como a medicina, agricultura, estudos da terra, entre outras. Uma área atualmente desafiadora e que desperta interesse para os pesquisadores, é a classificação de imagens hiper-espectrais [1]. A complexidade dos dados, em termos de bandas espectrais, torna difícil a sua classificação usando os algoritmos tradicionais, levando os pesquisadores a recorrerem a algoritmos avançados, como por exemplo, as redes neuronais.

As imagens de detecção remota hiper-espectrais dão-nos características muito específicas sobre a composição de uma certa região da terra. Comparando com uma imagem comum, estas oferecem informações mais ricas e detalhadas [2], daí o interesse no uso destas, para classificação e consequente mapeamento de zonas propensas a queimadas, ou mesmo a identificação de zonas queimadas.

O conhecimento do terreno, a sua composição e características, permite um melhor controlo do mesmo. Em [3], o autor afirma que, não é possível parar a natureza, mas pode-se minimizar a incidência de incêndios florestais e os danos que eles causam através do mapeamento das zonas de risco, para isso, é necessário, porém, identificar essas zonas.

Durante muito tempo, a obtenção de informações contidas numa imagem, foi feita ao alcance da visão humana. Nos dias de hoje, técnicas computacionais, como as de *machine learning*, são usadas e com resultados satisfatórios.

A principal motivação deste trabalho é de aplicar modelos/algoritmos de *machine learning* que, recebendo imagens de diferentes sensores, sejam elas arquivadas ou recebidas em tempo real, o utilizador final possa rapidamente obter informação/medidas sobre áreas

com riscos alto, médio e baixos em relação a ocorrência de incêndios, ou ainda, informação de zonas que mais foram atingidas por incêndios durante um dado período do ano.

## 1.2 - O Projeto

O projeto a que nos propomos, consistirá em realizar experiências na aplicação de quatro algoritmos diferentes de classificação sobre um conjunto de dados existente e comparar os seus resultados, tirando ilações dos mesmos. Num segundo momento do projeto, após os algoritmos serem treinados, seguir-se-á a fase de previsão das classes aprendidas na imagem real, ou seja, identificar partes da imagem que correspondam a áreas propensas a incêndios florestais. Essas áreas podem ser zonas com grande vegetação, zonas com proximidade a áreas de habitação populacional e ainda zonas propícias para a prática da agricultura.

Identificadas essas áreas na imagem, embora não seja objetivo principal deste projeto, fica em aberto a possibilidade de ser feito o mapeamento das zonas por forma a ter melhor controlo sobre elas, principalmente em períodos que se verifica maior incidência de incêndios, na maior parte em períodos secos.

A análise de imagens hiper-espectrais envolve conhecimentos em algumas áreas do saber que se relacionam entre si, entre elas, a análise de imagem e *machine learning*. A primeira vai ser objeto de estudo no capítulo de conceitos fundamentais, como forma de adquirir conhecimentos básicos, na obtenção de características e padrões de uma imagem e a segunda consistirá em aplicar algoritmos de classificação sobre um conjunto de dados.

## 1.3 - Objetivos

Com esta dissertação pretende-se aplicar modelos na classificação de imagens hiper-espectrais para identificação de áreas propensas a queimadas em Moçambique e posteriormente o seu devido mapeamento. Para materializar este objetivo, será necessário especificamente o seguinte:

- Obtenção de imagens anotadas;
- Análise das imagens, de acordo com técnicas e algoritmos apropriados;
- Classificação das imagens de acordo com as características obtidas no processo da análise.



## 1.4 - Área de estudo

Para condução deste trabalho, foi definida como área de estudo o distrito de Mandimba na província do Niassa - Moçambique. A província do Niassa, a mais extensa do país com cerca de 129 mil km<sup>2</sup>, localiza-se na região noroeste de Moçambique, entre as latitudes 11°25' norte e 15°26' sul e as longitudes 35°58' este e 34°30' oeste, faz fronteira a norte com a Tanzânia, a oeste com a República do Malawi, a leste com a província de Cabo Delgado e a sul com as províncias de Nampula e Zambézia [4]. A província do Niassa tem uma população estimada em 1.213.398 habitantes [5], e de acordo com o documento citado, é a menos populosa do país. Sendo a mais extensa e menos populosa, claramente que a maior parte dela é de cobertura vegetal, o que propicia a ocorrência de incêndios, tanto os provocados pela ação humana (com maior percentagem), que muitas vezes têm origem pelo desmatamento das áreas florestais para prática da agricultura e da caça, como os de origem natural. Por sua vez, o distrito de Mandimba localiza-se na zona austral desta província, com uma superfície de 4,699km<sup>2</sup> e uma população de 133,648, recenseada em 2007 [5]. O distrito faz fronteira com o distrito de Ngauma a norte, com os distritos de Majune, Maúa e Metarica a leste, com os distritos de Cuamba e Mecanhelas a sul e a oeste com a república do Malawi, uma fronteira de 110Km.



a)



b)

Figura 1-1 Mapa de Moçambique a), província do Niassa b) com destaque o distrito de Mandimba<sup>1</sup>

## 1.5 - Contribuições

Tendo em conta os objetivos desta dissertação, espera-se que os seus resultados venham a contribuir de maneira eficaz, em diferentes aspetos que envolvam a classificação de imagens hiper-espectrais, em particular no campo de deteção remota para diversos fins. Dessas contribuições podem-se destacar:

- uso de modelos computacionais, para uma rápida identificação das zonas propensas a ocorrência de incêndios;
- tendo as zonas identificadas, poderá ajudar no mapeamento das mesmas e, registo temporário dessas áreas;

<sup>1</sup> Imagem extraída da internet “divisão administrativa de moçambique e mapa da província do niassa” e destacada manualmente, com um editor de imagens, o distrito de Mandimba.

- apoio no controlo das áreas mapeadas e consideradas vulneráveis em termos de ocorrência de incêndios;
- apoio na tomada de decisão, perante novas situações que possam ocorrer de forma ocasional.

## **1.6 - Estrutura da dissertação**

Esta dissertação, na sua estrutura é composta por quatro capítulos, precedidos de uma introdução e que por fim são acompanhados de uma conclusão.

O capítulo 1, é reservado para o desenrolar da parte introdutória do trabalho, começando por um breve enquadramento da dissertação e descrição do projeto, apresentação dos objetivos, descrição da área de estudo e de algumas contribuições que o projeto pode suscitar.

O segundo capítulo apresenta alguns conceitos fundamentais, à volta de tratamento de imagens para aplicação em modelos de classificação, encontra-se ainda neste capítulo a revisão bibliográfica, onde são apresentados alguns estudos realizados sobre, riscos, monitoramento de incêndios e mapeamento de áreas consideradas de risco.

O terceiro capítulo foi reservado para apresentação da metodologia. É neste capítulo que se vai explicar sobre o estudo experimental, começando por apresentar uma resenha sobre a aquisição dos dados e a composição dos mesmos, a seguir pode-se ver os resultados alcançados para cada modelo proposto e a comparação entre os diferentes modelos.

Por fim, no quarto capítulo, sumaria a conclusão, apresentando os principais resultados, discussão sobre trabalhos futuros e apresentação das dificuldades encontradas no decorrer da dissertação.

## Capítulo 2

# Conceitos fundamentais e Revisão bibliográfica

### 2.1- Imagem hiper-espectral vs multiespectral

Durante muito tempo, os pesquisadores na área de detecção remota e outras que fazem o uso de imagens de satélite, baseavam-se em imagens multiespectrais; há poucas décadas, as atenções viraram-se para análise e obtenção de informações em imagens híper-espectrais [6].

A principal diferença entre imagens multiespectrais e hiper-espectrais está no número de bandas e o quão limitadas são as bandas [7]. As imagens multiespectrais geralmente têm entre 3 a 10 bandas representadas por pixels onde cada banda é adquirida com base num radiômetro de detecção remota. Por outro lado, as imagens hiper-espectrais consistem em bandas muito mais estreitas (normalmente de 10-20nm) e podem ter centenas ou milhares de bandas, adquiridas com um espectrómetro de imagem [6], [7].

A Imagem hiper-espectral faz parte de um conjunto de técnicas que, muitas vezes são referidas como análise ou imagem espectral [8], e, para gerar as imagens híper-espectrais é necessário passar pela fase de coleta e processamento de informações de todo o espectro eletromagnético (intervalo completo de todas as possíveis frequências da radiação eletromagnética).

As imagens hiper-espectrais têm aplicação em vários campos de estudo, tais como na exploração de minerais, na gestão de recursos e na monitoria ambiental [8]. São usadas em pesquisas na área da agricultura para monitorar o desenvolvimento e a saúde das culturas, na área militar, em pesquisas sobre vegetação, análise de alimentos, entre outras.

No que se refere ao espectro, e tendo como objetivo analisar ou estudar a vegetação, esta reflete mais na zona próxima de infravermelho e menos na zona de luz vermelha, comparativamente a estudos do solo [6].

O tamanho da área do solo representada por um único conjunto de medidas espectrais, define a resolução espacial da imagem e depende essencialmente do desenho do sensor e da sua altura em relação a superfície [6].

**Deteção remota** - a ciência e a arte de obter informação sobre uma área, um objeto ou um fenómeno por análise de dados que são adquiridos por um dispositivo que não esteja em contacto com o objeto em estudo [9].

## 2.2- Análise de Imagem

Como qualquer processo de análise, a análise de imagem tem como objetivo extrair características ou informação considerada importante, numa imagem. É importante sublinhar que, ao falar de imagem, nos referimos a imagens digitais, daí que para tal processo de análise seja necessário a aplicação de técnicas de processamento digital de imagem [10].

Embora a visão humana seja importante nesse processo, ela fica limitada quando há necessidade de analisar grandes quantidades de dados. O processamento digital de imagens não é uma atividade humana, ele apenas define os padrões que quer obter na imagem, mas é sempre indispensável a presença de computadores, tendo em conta que, algumas características são quase impossíveis de perceber pela visão humana.

### 2.2.1 Análise de Imagem Digital

Na análise de imagem digital, o principal elemento envolvido, é o computador. É uma verdade que, para melhor tratamento de informação/dados extraídos de uma imagem digital, o uso de sistemas computacionais nos trará medidas mais rápidas e precisas, e ainda garantem uma certa confiabilidade, tratando-se de medidas que manualmente são impossíveis de executar.

## 2.3- Processamento Digital de Imagem

Os autores [11] começam por definir a imagem como uma função bidimensional  $f(x,y)$ , com  $x$  e  $y$  sendo coordenadas espaciais em qualquer par de coordenadas  $f(x,y)$ , é chamada de *intensidade de imagem* nesse ponto. Quando os valores dessas coordenadas são finitos a imagem é chamada de *imagem digital* e a área que se responsabiliza em tratar desse tipo de imagem é o *processamento digital de imagem*.

O *processamento de imagem* trata-se de um processo em que tanto a entrada como a saída são imagens. Os autores acima referenciados consideram três tipos de processos

informatizados, sendo: processo de nível baixo, nível médio e de nível alto. De baixo nível, a entrada e saída são imagens; de nível médio, a entrada pode ser uma imagem mas, na saída temos atributos extraídos dessa imagem e, finalmente o de alto nível em que de um lado são conjunto de objetos reconhecidos (análise de imagem) e do outro lado são funções cognitivas, na sua maioria associadas a visão. Portanto, os autores definem *processamento digital de imagem* como processos que se pode ter entrada e saída, uma imagem, e ainda processos que extraem atributos de imagens e inclui o reconhecimento de objetos individuais.

Qualquer sistema de processamento de imagens, é composto por diferentes etapas [12], organizadas de forma sequencial, e destas destacam-se fundamentalmente cinco etapas, como ilustra a figura 2-1. A seguir são descritos o que cada etapa executa até o resultado final.

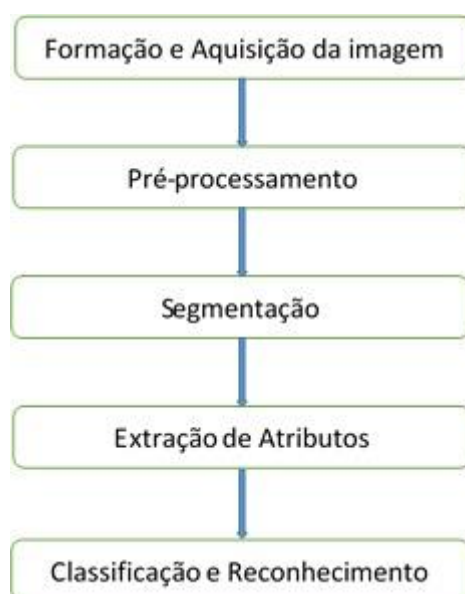


Figura 2-1 Principais etapas de um Sistema de Processamento Digital de Imagem, adaptado de [12]

### 2.3.1 Aquisição da Imagem

Na aquisição de imagens digitais, que também é chamado de detecção/sensoriamento, dois elementos principais são necessários [12], [13]: o primeiro a considerar é o dispositivo físico sensível ao espectro de energia irradiada pelo alvo que nos interessa, pode ser por exemplo o espectro de raio-x, luz infravermelha, luz ultravioleta, ou luz visível. É importante que o dispositivo na sua saída produza um sinal elétrico que seja proporcional ao nível de energia percebido; o segundo elemento é um dispositivo conversor, usualmente denomina-se por digitalizador, a sua principal função é converter o sinal elétrico analógico produzido pela saída do sensor para um formato digital.

### 2.3.2 Pré-processamento

A fase do pré-processamento consiste na melhoria da informação visual da imagem, de modo a tornar fácil a interpretação humana ou de máquina [12], [14].

Nas técnicas de pré-processamento deve-se considerar duas principais categorias, em que uma delas usa métodos que lidam com o domínio espacial e a outra usa métodos que lidam com o domínio da frequência. No primeiro caso, as técnicas baseiam-se em filtros que trabalhem com o plano da imagem e, no segundo caso, as técnicas de processamento baseiam-se em filtros que manipulam o espectro da imagem. Para dar maior destaque às características da imagem, a maneira mais usual é a combinação de diferentes métodos das duas categorias referidas. A imagem a seguir é basicamente um exemplo que ilustra um pré-processamento simples:



Figura 2-2 Aplicação do filtro de mediana para remoção de ruído [15]

Este é apenas um exemplo de vários filtros que se podem aplicar. No caso em que haja necessidade de dar mais realce ao contorno ou bordas de uma imagem, pode-se usar por exemplo o filtro passa-alta, filtro este que permite a passagem das frequências altas e impede ou reduz a passagem das frequências baixas, através da chamada frequência de corte [16].

### 2.3.3 Segmentação

Na etapa de segmentação é onde a imagem é dividida em regiões ou objetos que a constituem, ou seja, são extraídas partes dela que são essenciais para o processo de análise. Neste processo são usados métodos automáticos ou semiautomáticos para a extração das medidas, características ou informação da imagem que revele algum interesse na análise.

Vários trabalhos [11]-[13] consideram essa etapa como crucial, porque a obtenção de melhores resultados nas etapas seguintes e no final do processo, dependerão exatamente de como serão definidas as regiões de interesse para posterior processamento e análise, nesta etapa. Portanto, deve-se evitar ao máximo que nesta etapa se cometam erros, os mesmos podem contribuir para que se obtenham resultados não desejados e para uma ineficiência de todo o processamento.

Os algoritmos de segmentação de imagens, basicamente baseiam-se em uma de duas propriedades dos valores de intensidade: *descontinuidade* e *similaridade* dos níveis de cinza [11], [13]. A primeira propriedade preocupa-se com a partição de uma imagem, baseando-se em uma mudança brusca na intensidade, tal como acontece nas arestas de uma imagem. A segunda abordagem, tem a ver com a partição de uma imagem, em regiões similares, de acordo com um conjunto pré-definido.

Para as duas abordagens existem diferentes técnicas de segmentação que são usadas. Uma das técnicas amplamente usadas e que é baseada na similaridade é a binarização de imagem ou *image thresholding*, ao nível computacional é considerada simples e eficiente. Ela é adequada para situações em que as amplitudes dos níveis de cinza são suficientes para obter as características dos objetos presentes na imagem. Nesta técnica cada nível de cinza é o limiar de separação entre pixels que compõem os objetos e o fundo. O facto de ser chamada binarização, significa que na sua saída se obtém uma imagem binária, ou seja, uma imagem a preto e branco.

A técnica mais usada na abordagem de descontinuidade é chamada de detecção de bordas. Numa imagem de interesse as bordas caracterizam os contornos dos objetos que se encontram nela e cada ponto da borda é entendido como as variações bruscas dos níveis de cinza, contendo neles as características de transições entre objetos diferentes.

A figura 2.3 o processo de segmentação por binarização (B) e segmentação por detecção de bordas (C).

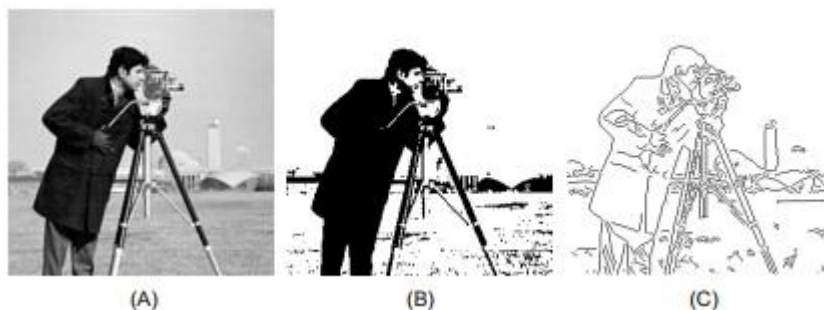


Figura 2-3 (A) Imagem original, (B) Segmentação por binarização e (C) Segmentação por detecção de bordas [17]



### 2.3.3.1 Morfologia Matemática

A morfologia é também uma das técnicas bastante aplicadas no processo da segmentação. O termo morfologia, é bem conhecido na área da biologia, sendo referido ao campo que estuda a forma e estrutura das plantas e animais. No contexto em causa (morfologia matemática), refere-se a uma ferramenta para extração de componentes de uma imagem, que poderão ser usados para representar e descrever a forma de certas regiões, como por exemplo, os limites [11], [18]. As técnicas da morfologia matemática são aplicadas para o pré-processamento e o pós-processamento, no primeiro caso pode ser por exemplo a filtração morfológica.

No processamento de imagens, para interpretação dos objetos presentes na imagem, a morfologia matemática socorre-se na teoria de conjuntos, onde, cada objeto representa um conjunto ou um subconjunto. Os elementos dos conjuntos neste caso específico, são coordenadas de pixels que representam os objetos ou outras características importantes numa imagem.

#### 2.3.3.1.1 Conceitos básicos sobre teoria de conjuntos

Sendo  $A$  um conjunto do plano cartesiano  $Z^2$ . Se  $a = (a_1, a_2)$  é um elemento de  $A$ , escreve-se  $a \in A$ , da mesma forma, se  $a$  não é elemento de  $A$ , escreve-se  $a \notin A$ . A figura a seguir mostra em resumo algumas notações sobre teoria de conjuntos.

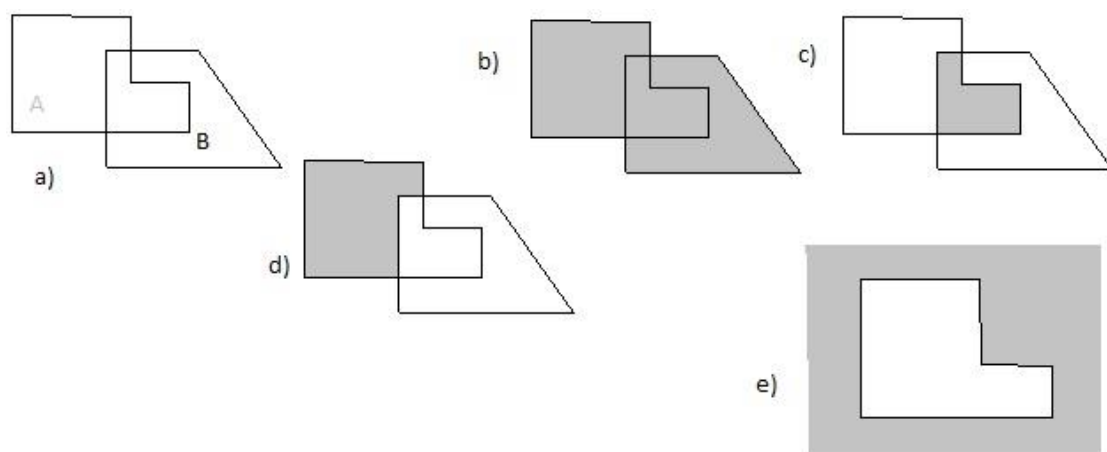


Figura 2-4 a) Representa dois conjuntos  $A$  e  $B$ . b) União entre  $A$  e  $B$ , denota-se por  $A \cup B$ . c) Intersecção de  $A$  e  $B$ , denota-se por  $A \cap B$ . d) Diferença de  $A$  e  $B$ , denota-se por  $A - B$ . e) Complemento de  $A$ , denota-se por  $(A)^c$ , figura adaptada de [11].

O complemento de um conjunto  $A$ , são todos os elementos que não estão contidos nesse conjunto:

$$A^c = \{w \mid w \notin A\} \quad (1)$$

A diferença entre dois conjuntos  $A$  e  $B$  ( $A - B$ ), é definido por:

$$A - B = \{w \mid w \in A, w \notin B\} = A \cap B^c \quad (2)$$

Na figura 2.4, os resultados das operações entre os conjuntos  $A$  e  $B$  é mostrado na cor cinzenta.

Na morfologia matemática, para além dos principais conceitos aplicados na teoria de conjuntos, adicionam-se a estes, mais dois conceitos: a reflexão e a translação [11].

A reflexão de um conjunto  $B$ , escreve-se  $\hat{B}$ , é definido como:

$$\hat{B} = \{w \mid w = -b, \text{ for } b \in B\} \quad (3)$$

A translação de um conjunto  $A$  por um ponto  $z = (z_1, z_2)$ , escreve-se  $(A)_z$ , é definida como:

$$(A)_z = \{c \mid c = a + z, \text{ for } a \in A\} \quad (4)$$

A figura 2.5, mostra os dois conceitos e o ponto a negrito é que identifica a origem dos conjuntos (as figuras usadas são as representadas na figura 2.4).

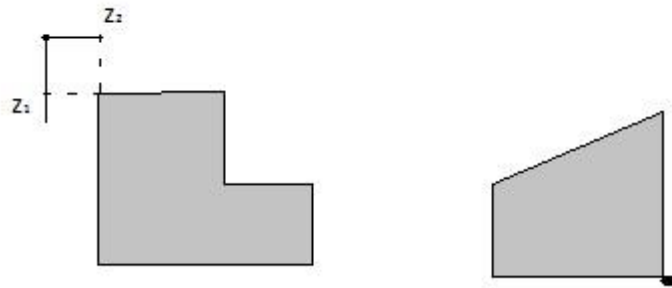


Figura 2-5 a) Translação de  $A$  por  $z$ . b) Reflexão de  $B$

Outro conceito que a morfologia matemática usa é o de operadores lógicos, envolvendo imagens binárias. Veja a tabela 2-1 que mostra os três principais operadores lógicos.

Tabela 2-1 Os três operadores lógicos básicos

$p$	$q$	$p \text{ AND } q (p.q)$	$p \text{ OR } q (p+q)$	$\text{NOT } p (\bar{p})$
0	0	0	0	1
0	1	0	1	1
1	0	0	1	0
1	1	1	1	0

Alguns operadores que também são usados e que, podem ser construídos a partir dos apresentados na tabela acima, são o XOR (OR exclusivo) e o NOT-AND. O XOR devolve 1 quando um dos pixels é 1, em casos de serem iguais, devolve 0; o NOT-AND (por exemplo  $[\text{NOT } (A)] \text{ AND } B$ ), seleciona os pixels que, simultaneamente estão em  $B$  e não em  $A$ .

Portanto, todas operações que foram citadas são aplicáveis na morfologia matemática, mas, existem duas outras operações, que são fundamentais e muito usadas, no processamento morfológico de imagens e com vários algoritmos que se baseiam nelas: a *erosão* e a *dilatação* [18].

**Erosão** - para que um elemento estruturante caiba numa imagem, a morfologia matemática considera fundamental a operação de erosão [11], [18], [19]. Para dois conjuntos  $A$  e  $B$  em  $Z^2$ , a erosão do conjunto  $A$  pelo conjunto  $B$ , denotado por  $A \ominus B$  e é definida através da seguinte equação

$$A \ominus B = \{z \mid (B)_z \subseteq A\}, \quad (5)$$

Onde  $\subseteq$ , representa a relação de inclusão de conjuntos.

**Dilatação** - sendo  $A$  e  $B$ , dois conjuntos em  $Z^2$ , a dilatação de  $A$  por  $B$ , denotado por  $A \oplus B$ , é definida através da seguinte equação:

$$A \oplus B = \{z \mid (\hat{B})_z \cap A \neq \emptyset\} \quad (6)$$

Esta equação é baseada na obtenção da reflexão de  $B$  pela sua origem e deslocando a reflexão por  $z$ . Logo, a dilatação de  $A$  por  $B$ , é o conjunto de todos os deslocamentos,  $z$ , dado que  $\hat{B}$  e  $A$  sobrepõem-se por pelo menos um elemento [11]. Desta interpretação resulta a seguinte equação:

$$A \oplus B = \{z \mid [(\hat{B})_z \cap A] \subseteq A\} \quad (7)$$

A dilatação é uma operação comutativa, ou seja,  $A \oplus B = B \oplus A$ :

$$\begin{aligned} A \oplus B &= \{z : z = a + b, a \in A, b \in B\} \\ &= \{z : z = b + a, a \in A, b \in B\} \\ &= B \oplus A. \end{aligned}$$

Em relação a complementação e reflexão de conjuntos, consideram-se, a dilatação e a erosão duais entre si:

$$(A \ominus B)^c = A^c \oplus B^c. \quad (8)$$

Ou seja,

$$(A \ominus B)^c = \{z \mid (B)_z \subseteq A\}^c. \quad (9)$$

## 2.3.4 Extração de Atributos

A extração de atributos é a etapa final de um sistema de processamento de imagem e que nos leva a classificação. É nesta etapa onde são extraídas as informações úteis da imagem processada [12]. Na parte intermédia existe a etapa chamada rotulação. Após a segmentação, a imagem é dividida em regiões correspondentes aos objetos e que são separadas das correspondentes ao fundo da mesma, neste ponto, as regiões são agrupadas uma próxima da outra por pixels que se tocam e, a seguir atribui-se um rótulo a cada grupo de pixel, o que permitirá para cada grupo de pixels próximos, calcular um parâmetro específico, pode ser a área, o perímetro ou outro que se considere adequado.

Por outro lado, temos os atributos da imagem, estes são divididos em duas classes de medidas: *características globais* - atributos da imagem toda, que pode ser por exemplo o número total de objetos e *características da região* - atributos de região, estes referem-se

aos objetos de forma independente, por exemplo a forma do objeto. Cada objeto identificado corresponde a um padrão e os valores medidos às características desse padrão. Conjunto de objetos semelhantes com características comuns, são pertencentes a mesma classe de padrões.

### 2.3.5 Classificação e Reconhecimento

A classificação é um procedimento responsável por classificar pixels e atribuí-los a categorias específicas, e, é ainda, considerada a tarefa principal e mais comum em *machine learning*, onde, o classificador é normalmente um mapeamento do tipo  $\hat{c}: \mathcal{X} \rightarrow C$ , onde  $C = \{C_1, C_2, \dots, C_k\}$  é um conjunto finito, geralmente pequeno, de rótulos da classe [20], e o uso do  $\hat{c}(x)$  significa estimação de uma função desconhecida  $c(x)$ . Durante o processo de classificação, são aplicadas diferentes técnicas, com fim de realizar o reconhecimento do objeto [13]. O objetivo do reconhecimento é na verdade realizar de forma automática a identificação dos objetos segmentados na imagem. No processo de classificação de formas são consideradas duas etapas principais, a de aprendizagem e a de reconhecimento. A figura 2.6, ilustra a representação geral para um processo de classificação.



Figura 2-6 Representação geral para um processo de classificação, adaptada de [13].

Consideram-se dois tipos de técnicas principais no reconhecimento de padrões: a classificação usando métodos de *aprendizagem supervisionada* e *não-supervisionada* e algoritmos da *aprendizagem supervisionada* são também divididos em duas categorias; os *paramétricos* e *não-paramétricos*. No caso da classificação com algoritmos *paramétricos*, o classificador é treinado com uma grande quantidade de amostra de padrões, cujas classes são conhecidas *a priori* [13], para estimação de parâmetros estatísticos como a média ou variância. Na classificação *não-paramétrica*, não são considerados os parâmetros estimados do conjunto de treinamento.

Na segunda técnica, classificação *não-supervisionada*, nada é definido *a priori*, apenas são criados *clusters* ou “agrupamentos naturais”, em português, dos padrões de entrada. Deve-se notar que diferentes tipos de algoritmos de *cluster* vão-nos dar diferentes *clusters* no

objeto. No caso de classificação de imagens multiespectrais, esta técnica agrupa os pixels de acordo com bandas espectrais que encontra [21].

No presente trabalho, serão usados e comparados os resultados de quatro algoritmos de classificação: *Random Forest Classifier*, *Gradient Boosting Classifier*, *Support Vector Machine* e *Deep Learning*, os quais são analisados com detalhes nas subsecções seguintes.

### 2.3.5.1 Random Forest Classifier

O classificador *Random Forest* (RF) - Floresta Aleatória, é uma combinação de árvores de decisão, de tal modo que cada árvore depende de uma amostra independente de valores de um vetor aleatório [22]-[24], tendo a mesma distribuição em todas as árvores na floresta. O RF considera a premissa de que um grupo de “aprendizes fracos” pode se unir e formar um “forte aprendiz” [25]. O erro generalizado na floresta de um classificador de árvores de decisão, depende da força das árvores individualmente na floresta e da sua correlação [22]. Melhorias importantes na acurácia da classificação resultam do crescimento no conjunto de árvores e fazendo com que elas votem na classe mais popular [22]. A classificação RF é um conjunto de classificações, isto é, refere-se a uma nova abordagem em que vários classificadores se juntam, ou seja, dezenas ou centenas de classificadores são construídos na classificação RF e as decisões são geralmente combinadas por um voto conjunto [24], por isto, a classificação RF é amplamente usada no processamento de imagens de detecção remota [24].

O classificador RF fornece probabilidade reduzida de ajustar variáveis explicativas dos dados de treino, ajustando de forma independente um grande número de árvores de decisão e o crescimento da árvore baseia-se no uso de subconjuntos aleatórios de dados de treino e um número limitado de variáveis de previsão selecionadas aleatoriamente [23].

### 2.3.5.2 Gradient Boosting Classifier

*Gradient Boosting* (GB), é um classificador que usa técnicas de *machine learning* para regressão e classificação e produz um modelo de previsão na forma de um conjunto de modelos de previsão fracos [26], um exemplo tipo, são as árvores de decisão [26], [27]. A ideia deste algoritmo veio do pensamento de que um aprendiz fraco, pode ser modificado e vir a ser melhor.

O objetivo do GB é de ensinar um modelo  $F$  para prever valores na forma  $\hat{y} = F(x)$  minimizando o erro médio quadrado  $(\hat{y} - y)^2$ , para valores verdadeiros  $y$ , valor este que é calculado sobre algum conjunto de dados de treino. Em cada fase do GB  $1 < m < M$ , pode-se assumir que existe algum modelo imperfeito  $F_m$ , e, o algoritmo GB não o muda de forma nenhuma, o que faz é melhorar o modelo adicionando um estimador  $h$  que forneça um melhor

modelo  $F_{m+1}(x) = F_m(x) + h(x)$ . O  $h$  pode ser calculado considerando que, de acordo com GB, se observa um  $h$  perfeito e, obtém-se  $F_{m+1}(x) = F_m(x) + h(x) = y$ , ou seja,  $h(x) = y - F_m(x)$ .

### 2.3.5.2.1 Algoritmo *gradient boosting*

**Entrada:** conjunto de treino  $\{(x_i, y_i)\}_{i=1}^n$ , uma função de perda diferenciável  $L(y, F(x))$ , número de iterações  $M$ .

Algoritmo:

1. Inicializa o modelo com um valor constante:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

2. For  $m=1$  to  $M$

- a. Calcular *pseudo-residuals*:

$$r_{im} = - \left[ \frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

- b. Ajustar um aprendiz base  $h_m(x)$  para *pseudo-residuals*, isto é, treiná-lo usando um conjunto de treino  $\{(x_i, r_{im})\}_{i=1}^n$ .

- c. Calcular o multiplicador  $\gamma_m$  resolvendo o seguinte problema de otimização unidimensional:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x) + \gamma h_m(x)).$$

- d. Atualizar o modelo:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

3. Saída:  $F_M(x)$ .

O algoritmo acima descrito, pode ser visto com mais detalhes em [26]-[29].

### 2.3.5.3 Deep Learning Classifier

São várias as definições que se atribuem ao algoritmo *deep learning* (DL), por exemplo, a que diz que DL é um ramo de *machine learning*, baseado em um conjunto de algoritmos que modelam níveis altos de abstração de dados [30]-[32], DL representa a era de crescimento de *machine learning*. O DL deve ser constituído de dois conjuntos de neurónios [32]: um que recebe o sinal de entrada e outro que envia o sinal de saída. Da entrada à saída, é composto por várias camadas, permitindo o algoritmo fazer o uso das múltiplas camadas de processamento. Quando a camada de entrada recebe uma entrada, ela passa uma versão modificada para a próxima camada, o que sugere que a saída de uma camada antecessora é a entrada da camada seguinte.

*Deep learning* como uma classe de algoritmos de *machine learning*, usa cascata de várias camadas de unidades de processamento não-linear para extração de características e sua transformação; é baseado na aprendizagem de múltiplas camadas de características e transformações de dados; é parte do mais amplo campo de *machine learning* na aprendizagem de representações de dados; aprendem vários níveis de representações que correspondem a diferentes níveis de abstrações. A composição de uma camada de unidades de processamento não-linear usadas num algoritmo de DL é dependente do problema a que se pretende resolver.

Os algoritmos de DL são baseados em representações distribuídas, isto é, dados observados são gerados por iterações de fatores organizados em camadas, e considera essas camadas como níveis de abstração ou composição [31], [32]. A figura a 2-7, ilustra uma breve arquitetura do algoritmo *deep learning*.

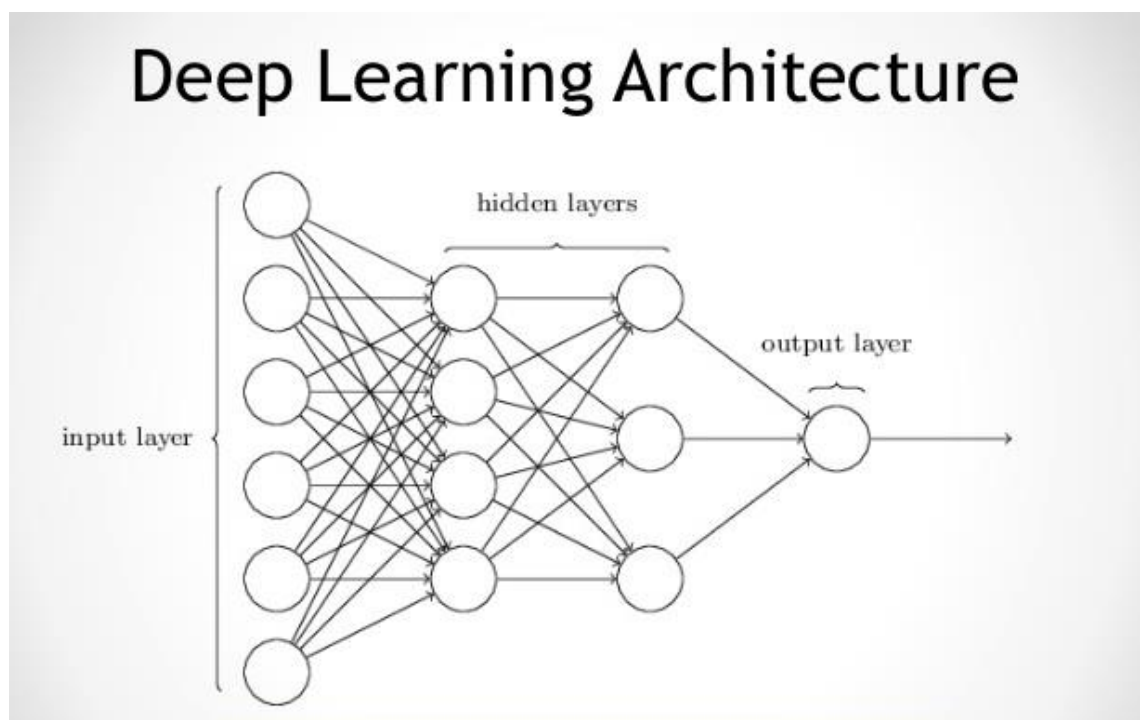


Figura 2-7 - Arquitetura deep learning [33]

O algoritmo de *deep learning* usa diferentes tipos de módulos. Neste trabalho, é proposto o módulo DBN (Deep Belief Networks), módulo este que consiste em dois diferentes tipos de redes neurais, BN (Belief Networks) e RBM (Restricted Boltzmann Machine) [34], que é considerado um modelo gráfico não direcionado de duas camadas e que apresenta um conjunto de unidades binárias ocultas [35]. Por outro lado as DBNs são uma composição de várias RBMs empilhadas uma por cima da outra [36]. O RBM em termos representativos é muito limitado, por isso, a necessidade de juntar vários deles para ter mais poder, dando origem desta forma o modelo DBN [35].

### 2.3.5.4 Máquina de Vector de Suporte

Outro algoritmo proposto para experiências neste trabalho é o SVM (Support Vector Machine), um algoritmo que foi durante muito tempo referência para classificação de imagens multi e híper-espectrais [37]. Este é um algoritmo supervisionado e poderoso de *machine learning* e é usado tanto para classificação como para regressão, embora o mais comum é a sua utilização em problemas de classificação [38].

Treinar um algoritmo SVM, consiste em encontrar a distância máxima dos padrões de treinamento mais próximos, ou seja, significa encontrar um híper-plano que separa com perfeição esses padrões [14]. Uma das vantagens desse algoritmo é que, mesmo que ele envolva otimização não linear no treino, a função objetivo é convexa e a solução de otimização é direta [39]. A figura a seguir ilustra a separação por híper-planos entre duas classes.

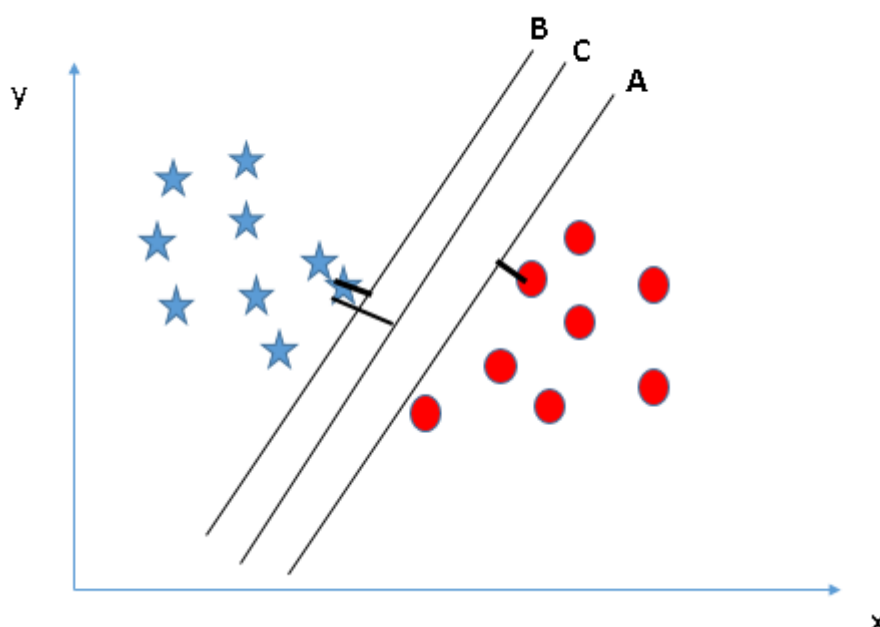


Figura 2-8 - Exemplo de SVM (separação por híper-planos) [38]

Observando a figura, é visível que a margem de separação dos híper-planos em relação as classes, o híper-plano C tem a melhor margem, por isso seria considerado o ótimo comparativamente aos A e B que, podem causar problemas de classificação errada.

Os SVMs baseiam-se em classificação binária [40], e quando se tem problemas de classificação com várias classes utiliza-se outras estratégias funcionais:

- ovo = One-Vs-One, ou seja, treina um classificador para cada duas classes individuais;
- ovr = One-Vs-Rest, ou seja, treina um classificador contra as outras classes em conjunto.

Neste trabalho em particular, trata-se de um problema multiclasse e uma das duas abordagens supracitadas será usada. Outra proposta para resolver problemas de multiclasse com algoritmos binários (no caso o SVM), é a abordagem ordinal, esta abordagem foi discutida



em [41], [42] com uso de hiper-planos ordinais e, consiste em pegar nas classes  $k$  que recebe e agrega em  $k-1$  classificadores binários, baseando-se na ordem das classes. Esta discussão pode-se encontrar em detalhes em [41] onde é proposto o método de replicação de dados. Este princípio de replicação de dados consiste em uma réplica dos dados originais para cada limite, se tomarmos em conta que as classes são separadas por hiper-planos.

## 2.4- Revisão Bibliográfica

O uso de imagens de satélites, para classificação e/ou mapeamento de áreas propensas a incêndios, é uma abordagem que despertou atenção para pesquisas há décadas. Um dos fatores que ajuda a obter bons resultados na classificação é a qualidade da imagem, onde se destacam as características das classes de interesse para o estudo [2], [43], [44]. [2], [43] defendem ainda que o risco de incêndio é bastante dependente do tipo de vegetação, o que é obvio afirmar que, os níveis de risco de haver ou não incêndios, pode dividir-se em alto, médio e baixo ou vários outros níveis dependendo do observador.

Atualmente, a classificação de imagens hiper-espectrais é amplamente usada em propostas de mapeamento de zonas com risco de ocorrência de incêndios [43]. Os autores em [43] realizaram um estudo sobre classificação da vegetação no mediterrâneo, mais concretamente na ilha de Elba (centro da Itália), para mapeamento das áreas com maiores riscos de incêndios, os autores propuseram a abordagem *Fuzzy*, por ser uma abordagem que melhor explora conjunto de dados integrados. Em termos de acurácia/exatidão (a acurácia varia dependendo do período da aquisição da imagem - primavera ou verão) na classificação, obtiveram um resultado satisfatório de 90.7% com dados adquiridos no período de verão. O modelo proposto peca por ter sido dependente de dados auxiliares; sem os dados auxiliares, obtiveram um resultado de 57.9% de acurácia, o que não representa uma boa performance do modelo proposto. Uma outra falha do modelo, foi que não mostrou especificamente a maioria dos tipos de vegetação do mediterrâneo, daí que, foi necessário adaptar de acordo com o conhecimento do terreno por parte dos autores. No presente trabalho pretende-se testar quatro algoritmos diferentes, anteriormente mencionados e encontrar o que melhor acurácia terá, utilizando primeiro um conjunto de dados (*dataset* de 2002), que serão divididos em treinamento e teste, e posteriormente os modelos serão testados com dados adquiridos em anos diferentes (*dataset* 2002 e 2005).

Para um melhor controlo de zonas que sejam propensas a riscos de incêndios, é fundamental que se tenha o registo das características sobre a cobertura florestal da área. O mapeamento da vegetação numa determinada área, pode servir de suporte para os decisores, no planeamento paisagístico, na modelação de risco de incêndios e posteriormente avaliar os impactos socioeconómicos de queimadas não controladas [45]. Dados de satélite, são os que

rapidamente podem providenciar informação de mapeamento e monitoramento florestais e, reduzir danos por incêndios e potenciais riscos [44]. Uma pesquisa desenvolvida pelos autores em [45] faz um estudo comparativo entre a classificação visual (visão humana) e classificação por máquina (visão computacional). Nesse estudo, em termos de acurácia, obtiveram melhor resultado na classificação por computador (no artigo não são apresentados dados percentuais sobre a comparação). Um dado importante nesse estudo é que, algumas características, por exemplo a de áreas queimadas, o modelo computacional que eles propõem não conseguiu obter informação real sobre a área mas, em contrapartida, por visão humana foi possível identificar os pontos que representavam essa característica no terreno. Portanto, os mesmos autores, propõem uma combinação entre a classificação por delimitação visual e a de máxima verossimilhança (por computador), para mapeamento de áreas com maior vegetação e propensa a incêndios na península de Halkidiki - Grécia (área onde foi desenvolvido o estudo), e, segundo o artigo, esta combinação na classificação, trouxe resultados com maior precisão para melhor manter o controle das zonas florestais com índices elevados de risco de incêndios, pelas instâncias competentes.

A visão humana é, sem dúvida, importante na obtenção de características constituintes de uma imagem, para posterior classificação, mas, atualmente, algoritmos poderosos e com resultados precisos são aplicados para diferentes tipos de classificação (textos, imagens, vídeos, entre outros). Neste trabalho, um dos algoritmos proposto, é o algoritmo de aprendizagem profunda (*deep learning*), ele é considerado bastante impressionante em termos de resultados, e tem sido cada vez mais aplicado por pesquisadores/estudiosos na área de aprendizagem computacional.

Os incêndios florestais, são a maior causa de degradação da terra nas regiões onde este fenômeno ocorre com maior frequência [43], [44]. Conhecendo as áreas e o período de maior incidência, é possível tomar medidas de prevenção, como a disponibilização de equipamentos de combate a incêndios e um alerta às entidades responsáveis, para que estejam preparadas para possíveis incidentes [3], [44].

Os autores [44], realizaram um estudo na zona Mediterrânea, região de Toscana - Itália, com objetivo de estimar e monitorar o risco de incêndios. Neste estudo foram usadas imagens do sensor NOAA-AVHRR NDVI - *Normalised Difference Vegetation Index*, para estimar o risco de incêndios naquela área. Valores de NDVI de associações de vegetação natural e seminatural, espera-se que sejam mais sensíveis indicadores a risco de queimadas do que os do tipo de cobertura agrícola.

Avaliados os tipos de vegetação derivados de mapa previamente existente, e, consideradas separadamente em novas análises de correlação, tiveram apenas melhorias não significantes e inconclusivas, daí que, os pesquisadores consideraram no seu modelo, uma estratégia supervisionada mais avançada, para identificar áreas onde as condições de vegetação derivadas do NDVI são mais relacionadas ao risco de incêndios. A estratégia

considerada foi a validação cruzada “leave-one-out”, e indicou que a precisão preditiva é realmente baixa nas resoluções espectrais mais altas, a baixa precisão em parte deveu-se à natureza da abordagem proposta. As estimativas de risco que se podem obter com os dados NDVI do sensor NOAA-AVHRR, dizem respeito as condições da vegetação, não tendo em consideração fatores importantes como as ações humanas, que podem causar incêndios após a aquisição de dados por satélite. Esses fatores são importantes ter em conta para uma boa precisão dos resultados.

Durante o período de verão em muitas partes do mundo, ocorrem incêndios florestais, uns induzidos pelo homem e outros por causas naturais. Recentemente, um estudo que tinha como objetivo a deteção de áreas propensas a incêndios, foi realizado em Kayer Khola, distrito de Chitwan, Nepal [3]. Os autores, na classificação de imagem de satélite, focalizaram-se na técnica de análise de imagem baseada em objetos, OBIA - *object-based image analysis*, com técnicas de modelação GIS - *Geographical Information System* e outros métodos de classificação, tendo em conta as informações espectrais, espaciais e de contexto, como também as propriedades hierárquicas. Para além dessas características, o OBIA fornece resultados precisos na identificação de florestas propensas a incêndios [3], e também permite melhorar a qualidade e a acurácia no processo de extração de características. O estudo revelou que 82% dos incêndios ocorrem em áreas florestais.

Áreas propensas a incêndios, são áreas que frequentemente são afetadas por chamas e que facilmente se espalham para outras áreas [3]. A deteção remota por satélite, abriu espaço para análise quantitativa de florestas e outros ecossistemas em toda a escala espacial e geográfica. Entender o comportamento dos incêndios florestais, os fatores que contribuem para um ambiente propenso a incêndios e fatores que influenciam no comportamento do incêndio, é fundamental para o mapeamento de áreas propensas a incêndios. A propensão a incêndios para qualquer área depende de vários fatores, entre eles, a cobertura da terra, precipitação, temperatura, topografia, proximidade a assentamentos e distância em relação a estradas.

Uma abordagem baseada em objeto, obtém melhores resultados de classificação, com um grau de acurácia alto, comparado por exemplo, com métodos baseados em pixel, a razão disto é o facto de o primeiro método explorar informação espectral e espacial, como se referiu anteriormente. Neste estudo sobre a floresta de Chitwan, a acurácia obtida na classificação foi de 86.54%.

Identificar focos na cobertura da terra, que são propensos a incêndios, tem uma grande importância para as entidades governamentais responsáveis pela gestão ambiental, não só pelo facto de ter essas áreas mapeadas, mas também para tirar medidas concretas sobre o quanto uma certa área foi devastada por queimadas num certo período.

Mapear essas áreas, é deveras importante, mas a questão fundamental é o tempo de vida útil do mapa, tendo em conta a ocupação de terras para prática da agricultura ou por

expansão urbana, o abate ilegal das florestas para exploração de madeira, entre outros fatores que são uma realidade para Moçambique [46].

A seguir é apresentada uma tabela resumo em relação aos resultados obtidos para cada caso de estudo acima discutidos.

Tabela 2-2 - Resumo de resultados sobre estudos passados

<b>Artigo</b>	<b>Metodologia/Algoritmo</b>	<b>Resultados (%)</b>
[3]	Object-based image analysis, técnicas de classificação e técnicas de GIS	86.54
[43]	Abordagem Fuzzy	90.70
[44]	Método baseado na identificação de pixels onde a concordância entre as variações inter- anuais na probabilidade de incêndios e os valores NDVI é o máximo possível. Foi também usada a Validação cruzada “leave-one-out” para validar os métodos comparando com estudos anteriores.	Resultado explicativo
[45]	Delimitação visual e máxima verossimilhança	Resultado explicativo

## Capítulo 3

# Metodologia & Resultados

### 3.1- Introdução

Para validação dos algoritmos propostos, serão consideradas as quatro modelações a seguir descritas:

**1 - Aquisição e análise dos dados** - neste primeiro ponto, o conjunto de dados disponibilizado, teve o tratamento necessário para utilização nos algoritmos de classificação, ou seja, as fases de pré-processamento e extração de atributos foram preparadas, estando disponíveis as classes de interesse para o estudo e os atributos necessários que, identificassem cada classe;

**2 - Treinar os modelos propostos** - neste ponto, com os dados fornecidos no ponto anterior, são treinados os modelos para reconhecimento das classes definidas para condução e com interesse no estudo;

**3 - Testar e avaliar os modelos** - após treinados os modelos no ponto 2, segue a fase do teste, em que, para este trabalho será usada a matriz de confusão (*confusion matrix*), para avaliar o quanto cada um dos modelos conseguiu aprender a reconhecer cada uma das classes consideradas;

**4 - Avaliar e analisar os resultados** - este é o ponto em que serão avaliados os modelos, tendo como métrica principal a acurácia (*accuracy\_score*), e, pretende-se ainda, analisar com que precisão os modelos conseguem prever dados novos, e desses resultados obter informação que nos permite tirar conclusões sobre os objetivos do trabalho.

### 3.2- Ferramentas utilizadas

Para a realização das experiências que são descritas nos subtemas a seguir, foi usada a linguagem de programação *Python* na sua versão 3.5, uma linguagem amplamente usada pela comunidade de *machine learning*, aproveitando as facilidades pelas bibliotecas e módulos que ela oferece. Podemos por exemplo, encontrar uma das mais completas bibliotecas para algoritmos de *machine learning*, a *sklearn*. No caso particular deste trabalho, optou-se por usar a plataforma *Anaconda* e o IDE (*Integrated Development Environment*) *Spyder* que, basicamente incorpora maior parte de bibliotecas e pacotes necessários para implementação ou aplicação de algoritmos de *machine learning*, tanto para classificação, como para regressão.

Quanto às características da máquina utilizada para testar os algoritmos, tratou-se de um computador com Sistema Operativo Windows 7 Professional Service Pack 1, 64 bits; Processador Intel core i7, CPU 2GHz; Memória RAM de 4G e uma placa gráfica NVIDIA.

### 3.3- Aquisição de dados

Para o presente estudo, foram usadas 7 imagens de uma mesma região, que representam 7 bandas espectrais diferentes. As imagens referidas representam a região do distrito de Mandimba, província do Niassa - Moçambique, adquiridas em 2002 e 2005, com uma cobertura da região de cerca de 283 Km<sup>2</sup>, ou seja, 571x551 pixels. A mesma região foi estudada por [47] e por [48]. A aquisição da imagem foi feita pelo satélite Landsat 7 (com resolução espacial de 30m), através da plataforma USGS *Earth Explorer* (<http://earthexplorer.usgs.gov>), nesta plataforma as imagens são adquiridas no formato TIFF - *Tagged Image File Format* com 16-bits, e com informação geoespacial necessária. Por razões de tamanho das imagens no formato TIFF, as imagens usadas neste trabalho foram convertidas para o formato PNG.

No processo de treinamento de um classificador, a primeira fase consiste na aquisição de dados. Neste estudo, para o treinamento dos classificadores propostos, serão usados um conjunto de dados que passaram por um processo de marcação supervisionada das características de interesse, onde um especialista fez a marcação do conjunto de pixels representando cada classe de interesse.

Neste trabalho foram usadas 5 bandas (verde, vermelho, próximo ao infravermelho, infravermelho de onda curta na faixa 1.55 - 1.75µm e infravermelho de onda curta na faixa 2.09 - 2.35µm) e duas composições, NDVI - *Normalized Difference Vegetation Index* e VI7 - *Vegetation index*, como mostra a tabela 3-1.

Tabela 3-1 - Descrição das bandas do satélite Landsat 7

ID	Banda Landsat 7	Comprimento de onda ( $\mu m$ )	Descrição <sup>2</sup>
1	2	0.52 - 0.60	<b>Verde</b> - enfatiza a vegetação de pico, que é útil para avaliar o vigor da planta
2	3	0.63 - 0.69	<b>Vermelho</b> - discrimina as encostas da vegetação
3	4	0.77 - 0.90	<b>Próximo ao infravermelho</b> - enfatiza o conteúdo de biomassa e linhas costeiras
4	5	1.55 - 1.75	<b>Infravermelho de onda curta</b> - discrimina o teor da umidade do solo e da vegetação e penetra em nuvens finas
5	7	2.09 - 2.35	<b>Infravermelho de onda curta</b> - alterações hidrotermais associadas a depósitos de minerais
6	NDVI	$NDVI = \frac{B_4 - B_3}{B_4 + B_3}$	Nesta composição das bandas vermelha e próximo ao infravermelho, resulta numa imagem que nos dá informação se a vegetação é verde ou não [48]
7	VI7	$VI7 = \frac{B_4 - B_7}{B_4 + B_7}$	Substitui o canal visível de vegetação clássica por canais infravermelhos de onda curta [48]

As imagens (em tons de cinza) a seguir, representam cada uma das bandas apresentadas na tabela 3-1.

<sup>2</sup> <https://landsat.usgs.gov/what-are-best-spectral-bands-use-my-study>

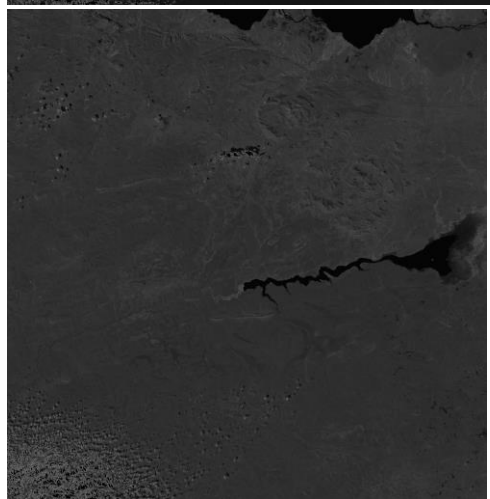
Tabela 3-2 - Imagens de cada banda espectral: a) banda 2, b) banda 3, c) banda 4, d) banda 5), e) banda 7, f) NDVI e g) VI7



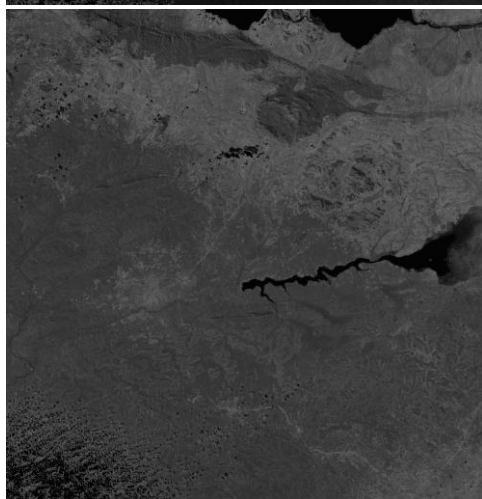
a)



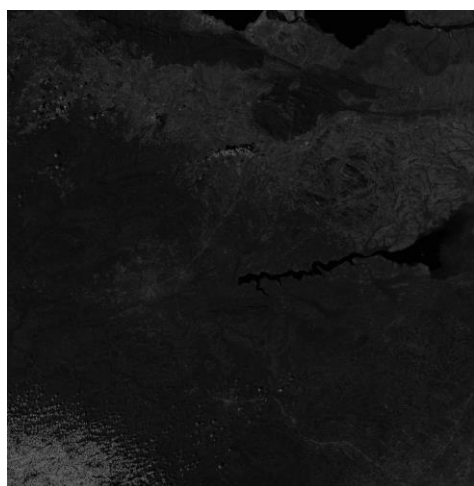
b)



c)



d)

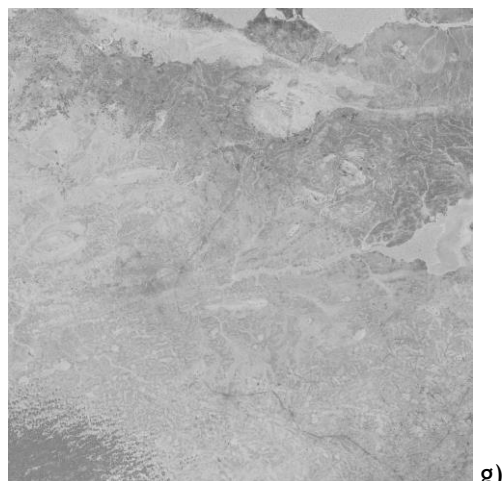


e)



f)





g)

As classes consideradas para este estudo, são as mesmas que foram usadas no estudo desenvolvido por [47], [48]. A tabela a seguir mostra uma breve descrição das classes.

Tabela 3-3 - Descrição das classes propostas

ID da classe	Nome da classe	Descrição
1	Corpos de Água	Áreas compostas apenas por água
2	Margens de Rios	Área próxima de mar/rio
3	Áreas Descobertas	Área seca, coberta por areia e vegetação seca
4	Terra de Cultivo	Área usada para prática da agricultura
5	Relvados/Pastagens	Grande área coberta de capim
6	Conjunto de Árvores	Quantidade densa de árvores e capim seco
7	Florestas	Grande área coberta de vegetação com árvores

Para cada classe, foram marcados um conjunto de pixels que correspondem à sua intensidade em cada uma das bandas espectrais, anteriormente mencionadas. No total são 7 características a observar (representadas por cada uma das bandas espectrais - Tabela 3-1) e um total de 9950 observações distribuídas conforme a tabela:

Tabela 3-4 - Número de observações por classe (dados 2002)

ID da classe	Número de observações
1	710
2	57
3	3428
4	309
5	714
6	339
7	4393

Ao conjunto de observações foi dividido em 70% para treino e 30% para teste em todos os classificadores. Nesta divisão foi usada a função em *python* `train_test_split(predictors, targets, test_size=.3)` que permitiu separar o conjunto de dados em treino e teste.

### 3.4- Experiência e resultados com modelos propostos

A experiência com os classificadores acima mencionados, mostrou que em termos de precisão eles foram muito próximos, embora seja visível uma ligeira diferença entre o DBN e os restantes, em todos os casos a primeira classe (corpos de água), foi bem classificada. O *Gradient Boosting* que ficou algumas décimas acima do *Random Forest*, teve má prestação no resultado, ao classificar as margens de rios/lagos (classe 2), com uma precisão de 36.84% e os restantes 63.16% ele classificou como áreas descobertas (classe 3 - área seca, coberta por areia e vegetação seca). A tabela seguinte ilustra a precisão/curácia a que cada classificador conseguiu alcançar.

Tabela 3-5 - Precisão/accuracy por classificador

Classificador	Precisão/curácia (%)
<i>Random Forest</i>	92.51
<i>Gradient Boosting</i>	92.63
<i>Deep Belief Networks</i>	<b>93.63</b>
<i>Support Vector Machine</i>	93.19

Analisou-se de seguida a precisão por classe, obtendo-se os resultados na tabela 3-6.

Tabela 3-6 - Precisão/accuracy do classificador por cada classe

Classificador	Precisão por classes						
	Corpos de Água	Margens de Rios	Áreas Descobertas	Terra de Cultivo	Relvado/Pastagens	Conjunto de Árvores	Florestas
<i>R. Forest</i>	100%	70.0%	99.0%	78.2%	70.2%	49.3%	94.6%
<i>G. Boosting</i>	100%	36.8%	98.0%	78.9%	70.8%	47.3%	96.4%
<i>DBN</i>	100%	83.3%	98.9%	79.2%	72.1%	42.6%	97.3%
<i>SVM</i>	100%	55.6%	98.6%	75.0%	70.5%	42.7%	96.4%

Os resultados apresentados na tabela 3-5, foram obtidos considerando os parâmetros por *default* (parâmetros pré-definidos) dos métodos usados. Estes resultados sugerem-nos que, os classificadores tiveram melhor comportamento para prever corpos de água (rios, lagos, lagoas, etc) todos a 100%; áreas descobertas previram mais de 98% e áreas com florestas (zonas cobertas de vegetação com árvores) previram entre 94% - 97.5%. Pode-se observar também que, para classes menos bem classificadas, a classe conjunto de árvores (quantidade densa de árvores e capim seco), foi a pior em todos os algoritmos, variando entre 42% - 49.5%.

### 3.5- Melhoria dos resultados

Para melhorar os resultados em algoritmos de classificação, várias técnicas podem ser aplicadas, uma delas, que o seu uso não necessita de muita experiência para aplica-la é o ajuste de parâmetros. O ajuste de parâmetros, é usado em algoritmos para controlar o seu comportamento [49], onde, cada valor usado no ajuste reflete um modelo diferente para esse algoritmo. Usou-se esta técnica para o *Random Forest* e para o *Gradient Boosting*, sem no entanto, serem observadas melhorias dignas de muita atenção.

Feito o ajuste de vários parâmetros, verificou-se que para o *Random Forest* os parâmetros `criterion='gini'`, `max_depth=50`, `class_weight='balanced'` e `n_estimators=50` e, para o *Gradient Boosting*, o parâmetro `learning_rate=0.2` com `n_estimators=200` influenciaram bastante na melhoria dos algoritmos embora com melhorias na casa de décimas. 30% dos dados foram usados para teste nos dois algoritmos. A tabela a seguir ilustra os resultados anteriores e os melhorados.

Tabela 3-7 Melhoria de RF e GB por ajuste de parâmetros

Classificador	Accuracy na forma clássica (%)	Accuracy na forma melhorada (%)
<i>Random Forest</i>	92.5	92.8
<i>Gradient Boosting</i>	92.6	92.9

### 3.5.1 Agregação de classes

Ainda como forma de melhorar os resultados dos classificadores, recorreu-se a agregação de classes, tendo sido considerada divisão por três níveis (baixo, médio e alto). Estes níveis representam o risco de ocorrência de incêndios. Nesta divisão, as classes foram separadas em  $\{1, 2, 3\} = 1$ ,  $\{4, 5\} = 2$  e  $\{6, 7\} = 3$ , seguindo a ordem crescente dos três níveis mencionados, como mostra a tabela:

Tabela 3-8 Agregação de classes em três níveis

Classes definidas no trabalho	Classes agregadas (propostas)
Corpos de água	1 - Baixo risco
Margens de rios/lagos/mares	
Áreas descobertas	
Terras de cultivo	2 - Risco médio
Pastagens	
Conjunto de árvores	3 - Alto risco
Florestas	

Na classificação por classes agregadas, foram considerados todos classificadores acima propostos e foram mantidos os ajustes de parâmetros realizados na subsecção anterior sobre o *Random Forest* e o *Gradient Boosting*. Os resultados não foram tão surpreendentes quanto o esperado, mas tiveram uma melhoria que vale a pena considerar, a acurácia teve uma ligeira subida de aproximadamente 2% para todos eles, continuando na frente o DBN com 95.71% e os restantes com aproximadamente 95% cada, como é possível visualizar na tabela que segue:

Tabela 3-9 Acurácia com classes agregadas

Classificador	Acurácia por classes			Acurácia geral
	1 - Baixo risco	2 - Risco médio	3 - Alto risco	
<b><i>R. Forest</i></b>	98.9%	74.5%	95.4%	94.9%
<b><i>G. Boosting</i></b>	98.4%	71.2%	97.1%	95.0%
<b>DBN</b>	99.0%	75.0%	97.0%	<b>95.7%</b>
<b>SVM</b>	99.0%	70.9%	96.8%	94.9%

Pelos resultados obtidos ao agregar as classes, as próximas experiências serão tomadas em consideração apenas três classes, sendo que, as mesmas compõem as classes de interesse para o estudo e as conclusões esperadas, ou seja, se se pretende prever áreas que sejam menos ou mais propensas a ocorrência de incêndios, as três classes propostas dão-nos informação de interesse, suficientes para as previsões.

### 3.5.2 Combinação de classificadores (*ensemble*)

O *ensemble* é também uma forma de melhorar resultados de acurácia nos algoritmos de classificação ou regressão, a sua abordagem está virada em pegar num conjunto de modelos aprendidos de forma individual e combiná-los para fazer previsões de novas instâncias [50]. Os métodos de *ensemble*, são algoritmos de aprendizagem, adaptados para construir conjunto de classificadores e, por meio de um voto ponderado das previsões, classificar novos pontos de previsão nos dados [50]. A principal ideia do *ensemble*, é que, não aprende com um simples classificador, mas sim, com um conjunto de classificadores e combina as previsões dos múltiplos classificadores.

Existem várias formas de construir os métodos *ensemble* para combinação de classificadores, neste trabalho, em particular, foi usado o método **VotingClassifier()**, disponível na biblioteca `sklearn.ensemble` da linguagem de programação *Python*. Para a experiência, consideraram-se os classificadores *Random Forest*, *Gradient Boosting* e *SVM*, e o resultado é o que se pode ver a seguir,

Tabela 3-10 Combinação dos classificadores RF, GB e SVM

Classificador	Acurácia por classes			Acurácia geral
	1 - Baixo risco	2 - Risco médio	3 - Alto risco	
<b>VotingClassifier</b>	99.1%	79.5%	95.2%	<b>95.1%</b>

Como mostra a tabela acima, comparando os resultados com os apresentados anteriormente, é possível perceber que a acurácia não se distanciou tanto da média entre os três algoritmos usados no *ensemble*, tendo se estabelecido a subida de -0.2%.

### 3.6- Uso de dados com diferença temporal (2002 vs 2005)

Neste subtema pretende-se, usando os mesmos algoritmos anteriores, realizar novas classificações, neste caso, considerando dados diferentes temporalmente. Os dados de 2002 serão usados para treino e os de 2005 para teste, sendo que os de 2002 tem um total de 9950 observações e os de 2005 com 10107 observações. A distribuição por classes não varia tanto, a segunda classe com menos observações e a primeira e terceira com observações quase equiparáveis. A seguir uma comparação por classes entre as observações dos dois anos em alusão.

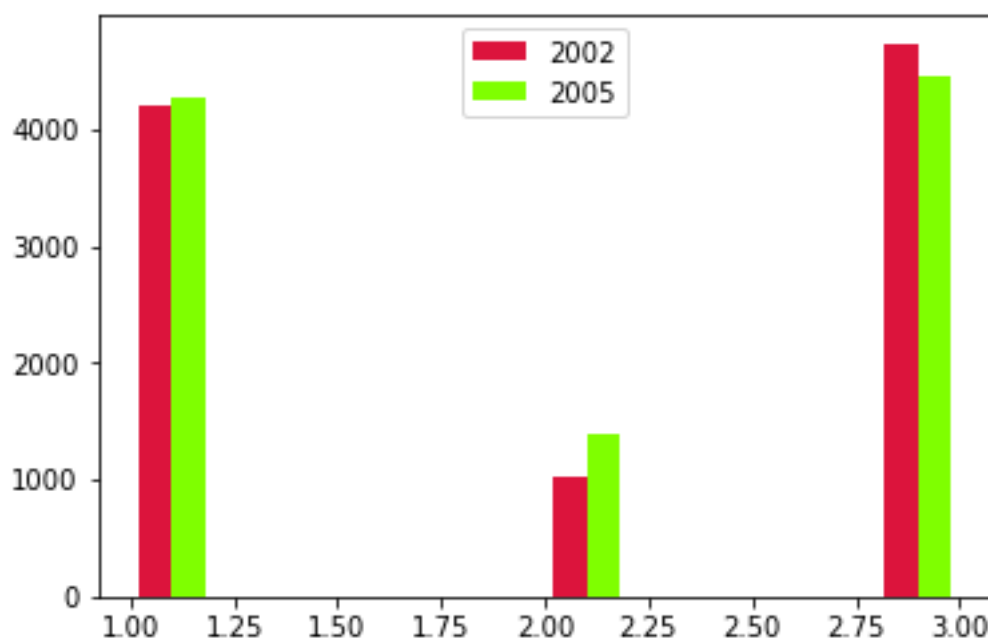


Figura 3-1 Visualização de observações por classes

A segunda classe (classe com risco médio) é a classe com menos observações nos dois conjuntos de dados, como se referiu anteriormente. Este facto, levou a má classificação para todos os algoritmos testados nesta classe. Os dados que a matriz de confusão nos mostrou, davam conta que, os classificadores ao classificarem mal a classe 2, parte considerável da previsão era considerada como pertencente a classe de maior risco (classe 3), do que propriamente a classe de risco médio (classe 2).

Nesta secção foram feitas experiências com os algoritmos acima citados. Surpreendentemente, em termos de resultados, o SVM e DBN, tiveram resultados muito próximos, deixando com resultados mais baixos o *random forest* e o *gradient boosting*, isto mostrou, mais uma vez o quão poderoso é o SVM. Para melhor perceção dos resultados, são apresentadas na tabela seguinte, a acurácia a que cada classificador conseguiu atingir, e através da matriz de confusão, os valores percentuais com que foram classificadas cada uma das classes como verdadeiras.

Tabela 3-11 Classificação com dados com diferença temporal (2002 vs 2005)

Classificador	Acurácia por classes			Acurácia geral
	1 - Baixo risco	2 - Risco médio	3 - Alto risco	
<i>R. Forest</i>	99.1%	50.8%	72.3%	80.7%
<i>G. Boosting</i>	99.6%	21.5%	77.1%	79.0%
DBN	99.9%	22.3%	86.5%	83.3%
SVM	99.4%	11.1%	91.9%	84.0%

### Dados 2005 vs 2002

A mesma experiência foi feita no sentido inverso, neste caso, foram considerados os dados de 2005 para treino e de 2002 para teste. Nesta experiência, enquanto que para os classificadores *random forest*, *gradient boosting* e SVM, os resultados foram melhores que na experiência anterior, tendo se verificado um salto de mais de 8% para o *gradient boosting* no resultado final da acurácia, o DBN teve o pior resultado ao sofrer uma queda de aproximadamente 5% comparativamente ao seu primeiro resultado. O facto de os dados de 2005 terem mais observações nas duas primeiras classes, pode ter influenciado bastante para obtenção de melhor acurácia nos primeiros três algoritmos referidos acima comparado com o primeiro caso. A tabela 3-11 ilustra de forma clara os resultados para cada um dos algoritmos.

Tabela 3-12 Classificação com dados com diferença temporal (2005 vs 2002)

Classificador	Acurácia por classes			Acurácia geral
	1 - Baixo risco	2 - Risco médio	3 - Alto risco	
<i>R. Forest</i>	95.8%	73.2%	70.3%	81.3%
<i>G. Boosting</i>	93.3%	26.3%	94.5%	87.0%
DBN	95.1%	75.0%	60.3%	78.6%
SVM	91.5%	45.8%	92.8%	87.4%

Como mostram as tabelas 3-10 e 3-11, o *random forest*, embora não tenha atingido a melhor acurácia, foi o classificador que, a nível das classes teve uma classificação mais equilibrada comparativamente aos outros três classificadores que, para a segunda classe estiveram abaixo de 50% em termos de precisão referente a primeira experiência, com exceção do DBN que, mesmo tendo sido o pior na acurácia geral no segundo caso, ele também mostrou um ligeiro equilíbrio em todas as classes.

Outra observação que mereceu atenção, é o facto de ter baixado a acurácia nos classificadores, se compararmos com os resultados obtidos com dados apenas de 2002, isto se verifica porque os modelos estão a tentar prever variáveis que provavelmente mudaram ao longo do tempo. Sabe-se que o mundo a cada minuto perde extensas áreas de florestas [51] por vários motivos (incêndios, expansão das zonas habitacionais, desmatamento para exploração de madeira, entre outros). Neste caso concreto, tendo características de imagens captadas em 2002 e em 2005, embora tenham sido captadas no mesmo período do ano [47], estas características podem variar ao longo do tempo, de região para região.

Um fenómeno por trás desta discussão (redução da acurácia verificada), é o chamado *concept drift*, discutido por vários pesquisadores. Este conceito diz que o facto de estar a fazer previsão com conjunto de dados que mudam ao longo do tempo, consequentemente sofre na sua acurácia, os resultados da previsão são menos precisos [52], [53]. Por outro lado, diz-

se que, em cada instante de tempo os dados de teste podem vir de distribuições diferentes dos de treino [53]. A figura abaixo ilustra uma simples arquitetura deste fenómeno

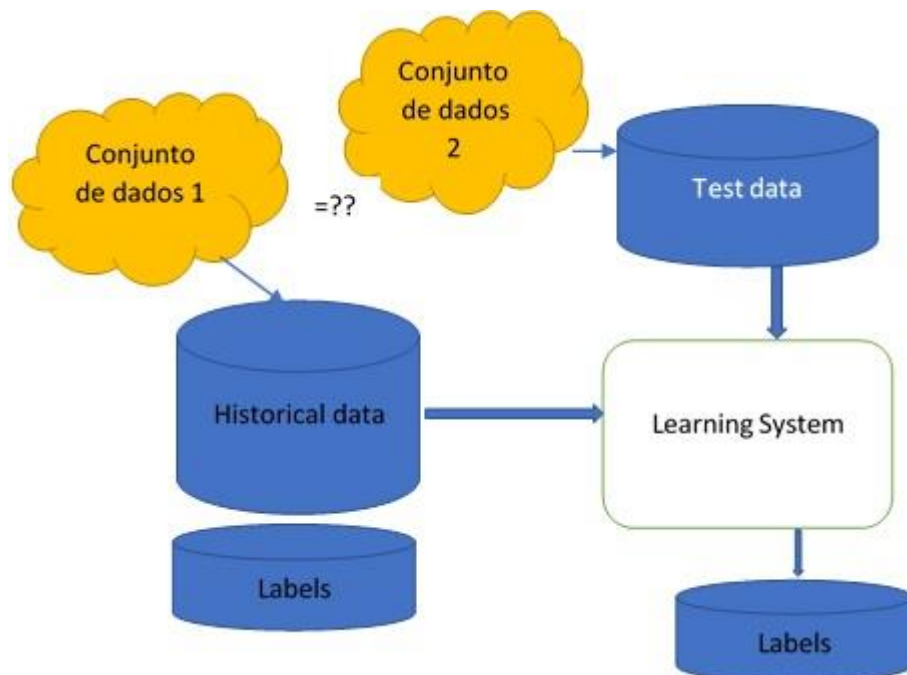


Figura 3-2 Arquitetura concept drift, adaptada de [53]

Contudo, apesar de verificar-se uma queda em relação a acurácia, pressupõe-se que, em termos de eficácia é bem melhor ao treinar e testar com dados de diferentes distribuições, do que com dados na mesma distribuição. Com dados na mesma distribuição, temos melhor acurácia, mas pode se dar o caso de o modelo estar a treinar e testar na mesma região ou na vizinhança.

### 3.6.1 Estratégia *ensemble*

Nesta experiência com conjunto de dados diferentes, também foi usada a estratégia *ensemble* como forma de melhorar os resultados da acurácia. Estes resultados não trouxeram melhorias, mas notou-se um equilíbrio em relação aos resultados por classificador. Nesta estratégia, em vez de usar um classificador típico do *ensemble* que fizesse o voto das classes melhor classificadas por cada classificador (ex.: `VotingClassifier()`), foi usada uma técnica em que o modelo que esteve melhor na acurácia, faz a classificação sobre os resultados dos outros modelos. Para isso, tendo em conta os resultados obtido acima, fez-se duas experiências, primeiro com o SVM como o modelo principal a seguir com o modelo DBN como o principal e, os resultados foram similares. Eis a seguir a tabela que ilustra o resultado obtido considerando o DBN como o modelo principal:



Tabela 3-13 Resultado do ensemble com dados diferentes no tempo

Classificador	Acurácia por classes			Acurácia geral
	1 - Baixo risco	2 - Risco médio	3 - Alto risco	
DBN ( <i>ensemble</i> )	Dados 2002 vs 2005			
	100%	29%	86%	84%
	Dados 2005 vs 2002			
	95%	76%	48%	71%

Os resultados apresentados na tabela 3-12, sugerem-nos que ao treinar com dados de 2005 e prever com dados de 2002 o DBN comporta-se melhor ao classificar as classes com mais observações em 2005 (classes 1 e 2) tanto de forma independente (tabela 3-11) como em *ensemble*.

### 3.6.2 Comparação dos algoritmos GB, RF e SVM

Para verificar o comportamento dos algoritmos, foi feita uma breve comparação entre os algoritmos GB, RF e SVM, onde é garantido que os mesmos estão a ser avaliados da mesma maneira [54]. Nesta comparação são considerados os valores da acurácia média e do desvio padrão apresentados a seguir:

Tabela 3-14 Comparação por acurácia média dos algoritmos GB, RF e SVM

Algoritmo	Acurácia média	Desvio Padrão
GB	0.81	0.26
RF	0.81	0.26
SVM	0.81	0.26

A seguir é ilustrada uma representação gráfica em relação a variação de cada um dos três algoritmos na tabela 3-13.

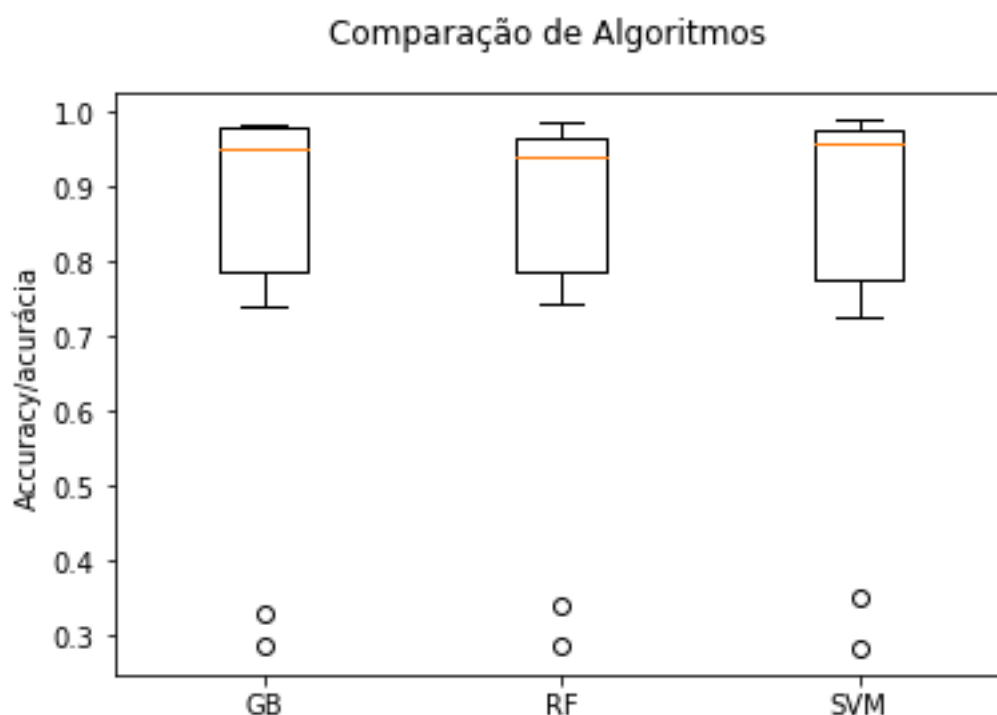


Figura 3-3 Comparação por acurácia média dos algoritmos GB, RF e SVM

Em termos comparativos, a tabela 3-13 apresenta valores aproximados da acurácia média e do desvio padrão. Como os valores são todos iguais não conseguimos tirar alguma informação para decidir qual dos três algoritmos teve melhor desempenho. Portanto, com o gráfico apresentado na figura 3-3 é possível notar que o SVM mostrou-se melhor com média ligeiramente superior aos dois outros modelos, e como todos têm o mesmo desvio padrão, concluímos que para todos possíveis valores da acurácia média que o SVM pode obter, ficam sempre próximos à média, este dado, dá-nos uma certa confiança ao considerar este modelo para classificação.

Outra conclusão que a figura sugere é que, em termos gerais o desempenho dos modelos escolhidos é semelhante, inclusive houve dois *folds* com pior desempenho em todos os modelos.

### 3.6.3 Visualização da previsão das classes na imagem

A informação que as métricas nos oferecem são de extrema importância no que se refere a análise e/ou classificação de imagens, mas, é também facto a necessidade de visualização das previsões dos modelos na imagem real. Nesta subsecção, são apresentadas as previsões na imagem, que, mostram o comportamento dos modelos *random forest*, *gradient boosting* e SVM, em relação as três classes previstas (baixo risco, risco médio e alto risco).

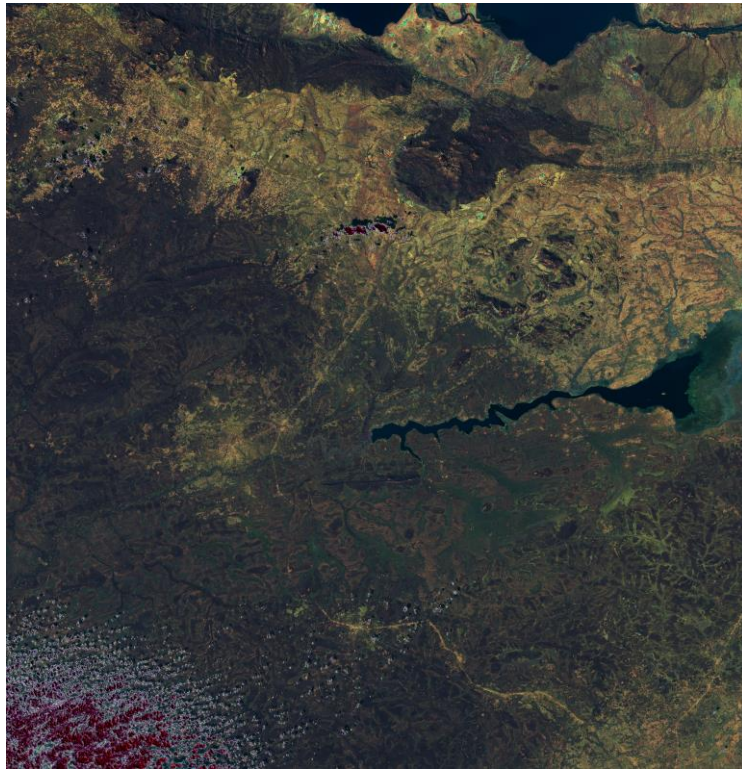


Figura 3-4 Imagem original representando a área em estudo, sem remoção de nuvens

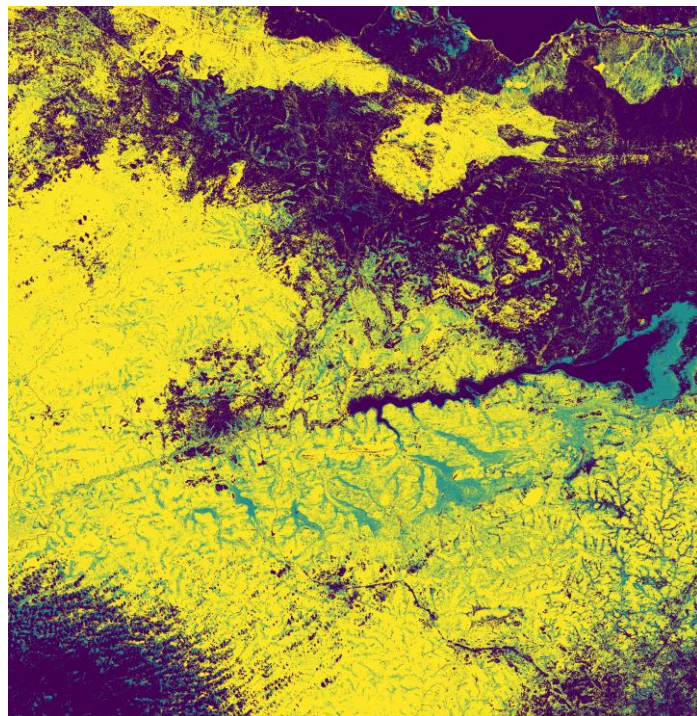


Figura 3-5 Previsões das classes com *Random Forest*



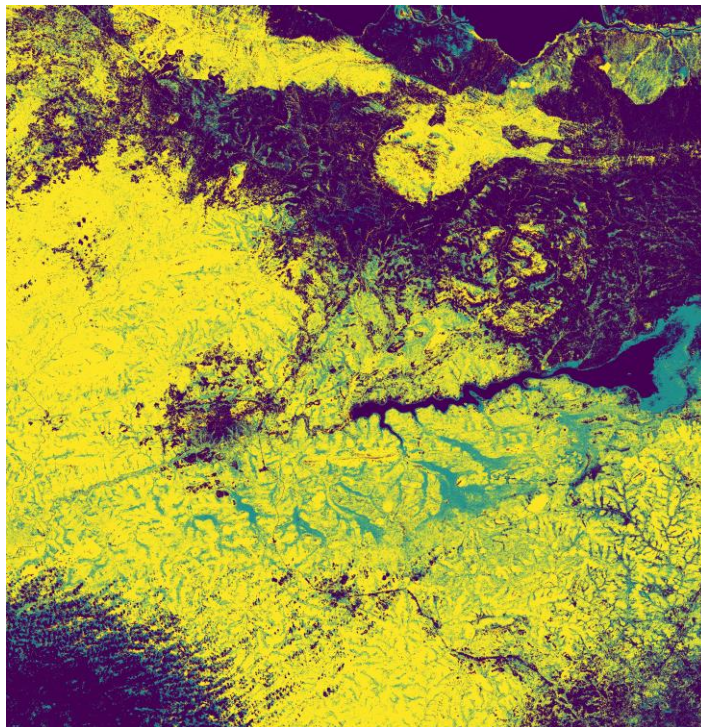


Figura 3-6 Previsões das classes com *Gradient Boosting*

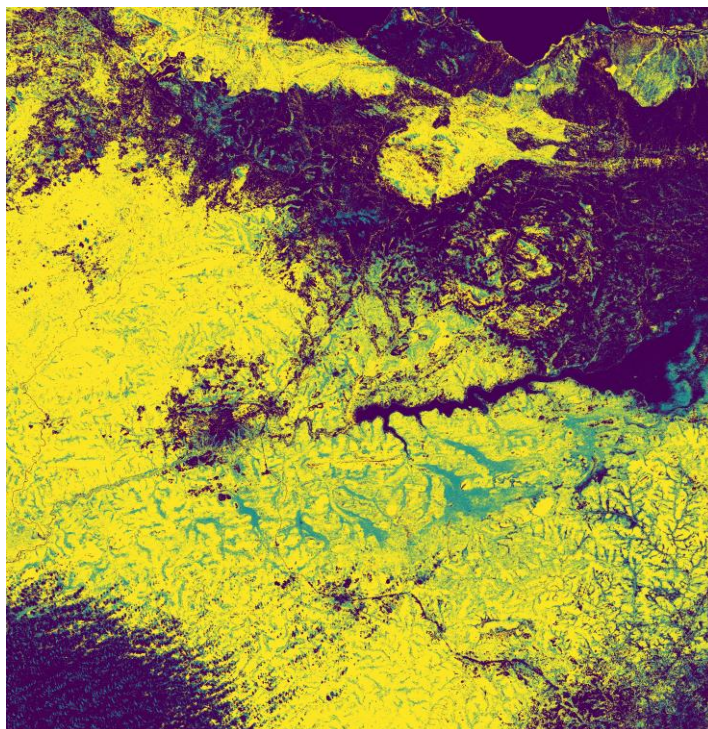


Figura 3-7 Previsões das classes com SVM

As três cores que são possíveis de visualizar nas imagens, representam as classes que se estão a prever, sendo, a cor roxa representa zonas com menor risco em termos de ocorrência de incêndios, a cor esverdeada representa zonas com risco médio e por fim o amarelo são as

zonas com maior vegetação e que se prevê maior risco no que se refere a ocorrência de incêndios.

Pode se observar nas imagens, que, o SVM perde algumas áreas em que os outros classificadores consideraram como de risco médio, ele considera como sendo áreas de menor risco, por exemplo, as margens de rios.

### 3.7- Entendendo o Amazonas pelo espaço

Esta experiência, baseou-se numa competição aberta pelo sítio da *kaggle*<sup>3</sup>, sítio este que, abre desafios, dando oportunidades para àqueles que desejam desenvolver competências no campo de *machine learning*, e, aos que já fazem parte desta grande comunidade, é uma oportunidade para mostrarem o seu potencial e criatividade e talvez conseguir alcançar o prémio que muitas vezes tem sido aliciante.

A competição tinha como objetivo usar imagens de satélite para rastrear a pegada (caminho) humana na floresta amazónica, tendo como finalidade apoiar o governo e outras entidades interessadas, a controlar de forma rápida, a invasão humana na floresta, que tem sido uma das principais causas do desmatamento naquela região.

Os dados fornecidos foram divididos em 40.479 imagens para treino e 61.191 para teste. De um lado tinham imagens num formato comprimido (.jpg) e do outro lado, imagens no formato de aquisição (.tiff) com quatro bandas (*red, green, blue & near infrared*), nesta experiência foram consideradas apenas as imagens no formato jpg.

Para resolução do problema, era necessário numa primeira fase, fazer a extração de atributos, tanto para as imagens de treino, como para as de teste, tendo sido implementado um método para obtenção de resultados nesta primeira fase do processo.

Após o tratamento dos dados referido no parágrafo acima, seguiu-se a fase de treinamento do modelo e realização das previsões de acordo com as classes que foram propostas para esta competição:

- *cloudy*
- *partly cloudy*
- *hazy*
- *primary rain forest*
- *water (rivers & lakes)*
- *habitation*
- *agriculture*
- *road*

---

<sup>3</sup> <https://www.kaggle.com/c/planet-understanding-the-amazon-from-space>

- *cultivation*
- *bare ground*
- *slash and burn*
- *selective logging*
- *blooming*
- *conventional mining*
- *“artisinal” mining*
- *blow down*

Para fazer a previsão dos rótulos na imagem, foi disponibilizado um ficheiro .csv com o nome da imagem e mesmo rótulo para todas as imagens, veja a seguir as primeira dez linhas retiradas do ficheiro referido.

```
image_name,tags
test_0,primary clear agriculture road water
test_1,primary clear agriculture road water
test_2,primary clear agriculture road water
test_3,primary clear agriculture road water
test_4,primary clear agriculture road water
test_5,primary clear agriculture road water
test_6,primary clear agriculture road water
test_7,primary clear agriculture road water
test_8,primary clear agriculture road water
test_9,primary clear agriculture road water
test_10,primary clear agriculture road water
```

Contudo, o resultado esperado para submissão, deve mostrar o nome da imagem e somente os rótulos que cada imagem contém. O resultado das primeiras dez linhas do ficheiro gerado para submissão foi o seguinte:

```
image_name,tags
test_0,primary clear
test_1,primary clear
test_2,primary partly_cloudy agriculture water
test_3,primary cultivation clear agriculture
test_4,primary partly_cloudy cloudy water
test_5,primary clear
test_6,primary partly_cloudy road agriculture habitation
test_7,primary road clear agriculture habitation
test_8,primary clear
test_9,primary haze cultivation clear agriculture
```

O modelo utilizado para treino e as respectivas previsões, foi o *XGBClassifier()*, esta função é importada de uma biblioteca otimizada da distribuição do *gradient boosting*, tecnicamente considerada eficiente e flexível chamada *xgboost*, forma abreviada de “*Extreme Gradient Boosting*”, como o próprio nome nos sugere, o *xgboost* tem como modelo de base o *gradient boosting* descrito numa das subsecções anteriores. Em relação ao *score* alcançado, o modelo não foi muito além de 88.2%. A figura 3-8 representa uma imagem real e a respetiva descrição sobre o que o modelo conseguiu prever nela.



Figura 3-8 Imagem de teste (test\_7) sobre uma parte do Amazonas. Previsões: primary, road, clear, agriculture, habitation

Como se pode observar na imagem, pelos rótulos que o modelo conseguiu prever é possível serem visualizados na figura 3-8, embora a visibilidade da imagem não seja clara, rapidamente conseguimos identificar rótulos como *habitation*, *road*, *agriculture*, os outros rótulos não são visíveis logo a primeira vista, mas com um pouco de atenção, localizam-se pequenos focos onde encontram-se representados.

# Capítulo 4

## Conclusão

Na introdução deste trabalho, foi previamente discutido o seu objetivo, que cingia concretamente em usar/aplicar modelos de *machine learning* na classificação de imagens hiper-espectrais para identificação de zonas consideradas vulneráveis em termos de ocorrência de incêndios. Como forma de alcançar este objetivo, foi definida a área de estudo que, baseou-se no *dataset* disponibilizado. A princípio podia-se optar por qualquer região de Moçambique, mas durante o processo de aquisição/preparação dos dados, foram fornecidos dados utilizados num outro estudo desenvolvido por [47] e [48]. A seleção dessa área de estudo, esteve em causa a sua densidade populacional (densidade populacional maior no sul da província do que no norte, o distrito de Mandimba localiza-se a sul) e o facto de ser uma área não coberta por nebulosidade [47].

Após a aquisição dos dados, seguiu-se a análise dos mesmos, tendo como foco as características de interesse na imagem. As características de interesse para o estudo foram extraídas em imagens coletadas em anos diferentes (2002 e 2005), e foram seleccionadas 7 classes distintas, mas iguais nos dois anos. Para o nosso objetivo, as classes foram reduzidas para três, agrupando as em menor risco de ocorrência de incêndios, risco médio e alto risco.

No decorrer das experiências, verificou-se que com classes agregadas obteve-se melhores resultados em termos de classificação, ou seja, os modelos testados estiveram melhor ao fazer a classificação com três classes do que com sete classes, mas como afirmado anteriormente, a agregação em três classes deveu-se aos objetivos do estudo.

Ignorando as questões de características/classes de interesses, esses resultados mostraram que, para classificação de *datasets* com várias classes, precisa de modelos mais avançados (por exemplo, os modelos de *deep learning*) ou de uma boa técnica de afinamento de parâmetros dos modelos conhecidos como tradicionais.

Embora os modelos “tradicionais”, como os que foram usados nesta dissertação (RF, GB e SVM), são amplamente usados na classificação de imagens multi ou hiper-espectrais, e obtendo bons resultados, os modelos avançados (redes neuronais, *deep learning*) são recomendados para tratamento deste tipo de dados, visto que, principalmente as imagens hiper-espectrais apresentam uma grande complexidade em termos de bandas espectrais.

As fontes de aquisição de dados para treinar e testar um modelo, são fatores que influenciam de alguma maneira no resultado da classificação. Nesta dissertação foi feita uma experiência com dados de 2002 para treino e de 2005 para teste, sobre a mesma região e coletados no mesmo período do ano, os resultados já eram de esperar, todos os modelos



tiveram uma queda na sua acurácia de pouco mais de 10% em relação a acurácia obtida apenas com dados de uma mesma fonte (2002).

No que se refere a previsão das classes na imagem real, foram utilizados três modelos para o efeito (RF, GB e SVM), e, analisando os resultados, verificou-se que o RF e GB, tiveram um comportamento similar na identificação das três classes previstas enquanto que para o SVM, algumas áreas em que para os dois primeiros eram de risco médio, ele considerou como sendo de menor risco (proximidade aos rios).

Pelos resultados referidos nos parágrafos anteriores, podemos afirmar que os objetivos previstos foram alcançados. Contudo, qualquer trabalho de natureza científica, deixa algumas reservas e interesses, na continuação da pesquisa como forma de melhor entender o/s assunto/s que foram aqui, objeto de discussão.

- a) **Trabalhos futuros** - no que diz respeito a estudos futuros, um particular interesse sobre classificação de imagens hiper-espectrais, passa por focar-se num estudo de raiz, desde a coleta de imagens e o seu devido tratamento, estamos aqui a falar da preparação dos dados para a sua utilização em diferentes modelos de classificação. Ainda na senda de estudos futuros, tendo a marcação de zonas suscetíveis a ocorrência de incêndios, o desafio será, mapear essas zonas ou ter registos em termos numéricos, dado que, os estudos atuais sobre incêndios em Moçambique, têm como referência sobre áreas queimadas, registos de 1990 [55], e claro, é necessário atualizar esses registos, visto que, ao longo dos anos mais áreas são queimadas/desmatadas.

Neste conjunto de desafios futuros, junta-se o estudo aprofundado de modelos de classificação avançados, como as redes neuronais, *deep learning* e os modelos nele envolvidos, tais como, as redes neuronais recorrentes, redes neuronais convolucionais, entre outros. Estes modelos, tem sido objeto de estudo e aplicação dos mesmos, por vários pesquisadores na área de análise e classificação de imagens e em outras vertentes da aprendizagem computacional.

- b) **Limitações** - para além das dificuldades decorrentes da pouca experiência neste campo de estudo, uma das principais limitações encontradas durante o percurso da dissertação, foi relacionada com os dados que foram disponibilizados. Como os dados foram pré tratados por outro grupo de pesquisadores, em algum momento sentiu-se a falta de informação detalhada destes, por exemplo, quanto a normalização, não nos foi dada informação suficiente que ajudasse a tratar de novos dados seguindo a mesma normalização. Estas limitações servirão como desafios para estudos futuros, por isso a necessidade de continuar com o estudo, começando com dados não processados, como se fez alusão em trabalhos futuros.

# Referências

- [1] B. Luo, J. Chanussot, and H. Blanche, "HYPERSPPECTRAL IMAGE CLASSIFICATION BASED ON SPECTRAL AND GEOMETRICAL FEATURES," 2009.
- [2] H. Li, Z. Ye, and G. Xiao, "Hyperspectral Image Classification Using Spectral - Spatial Composite Kernels Discriminant Analysis," vol. 8, no. 6, pp. 2341-2350, 2015.
- [3] K. Uddin, "Detecting Forest Fire Prone Areas Using Object-Based Image Analysis and GIS Techniques: A Case Study in Kayer Khola, Nepal," vol. 1, pp. 748-755, 2012.
- [4] Ministério da Agricultura - Moçambique', "Relatório de trabalho de campo realizado no âmbito do cumprimento das decisões de S. Excia. o Senhor Primeiro Ministro na sua visita à Província do Niassa." DNTF & IIAM, Lichinga, p. 60, 2010.
- [5] I. Instituto Nacional de Estatística, "Recenseamento Geral Da População E Habitação 2007 Indicadores Socio-Demográficos Da Provincia De Cabo Delgado," 2007.
- [6] W. Contributors, "Hyperspectral imaging," *Wikipedia, Free Encycl.*, pp. 1-24, 2015.
- [7] "Multispectral vs Hyperspectral Imagery Explained - GIS Geography," 2017. [Online]. Disponível: <http://gisgeography.com/multispectral-vs-hyperspectral-imagery-explained/>. [Acessado: 05-Apr-2017].
- [8] "Hyperspectral Imaging - Methods, Benefits and Applications," 2013. [Online]. Disponível: <http://www.azom.com/article.aspx?ArticleID=8495>. [Acessado: 05-Apr-2017].
- [9] C. Blodgett, "What is Hyperspectral Imagery ( HSI )? What is Remote Sensing ?, " *Missouri Resour. Assess. Partnersh.*
- [10] "Image Analysis," 2013. [Online]. Disponível: [https://en.wikipedia.org/wiki/Image\\_analysis](https://en.wikipedia.org/wiki/Image_analysis). [Acessado: 01-Nov-2016].
- [11] R. C. & Gonzalez and R. E. Woods, *Digital Image Processing*, Second Edi. Upper Saddle River: Prentice Hall: Tom Robbins, 2002.
- [12] O. M. Filho and H. V. Neto, "Processamento Digital de Imagens," *Rev. Bras. Geofísica*, vol. 21/03, no. 1, p. 331, 1999.
- [13] J. E. R. de Queiroz and H. M. Gomes, "Introdução ao Processamento Digital de Imagens," *Rita*, vol. 8, no. 1, pp. 1-31, 2001.
- [14] R. O. Duda, P. E. Hart, and D. G. Stork, *Pattern Classification*, 2nd ed. New York, 2001.
- [15] "Fundamentos de Computação Gráfica." [Online]. Disponível: <https://webserver2.tecgraf.puc-rio.br/~mgattass/fcg/trb12/ElianaGoldner/images.html>. [Acessado: 18-Nov-2016].
- [16] "Filtro passa-altas," 2013. [Online]. Disponível: [https://pt.wikipedia.org/wiki/Filtro\\_passa-altas](https://pt.wikipedia.org/wiki/Filtro_passa-altas). [Acessado: 18-Nov-2016].
- [17] "segmentação por detecção de bordas e por binarização - Google Search." [Online]. Disponível: [https://www.google.pt/search?q=segmentação+por+detecção+de+bordas+e+por+binarização&biw=1366&bih=662&source=lnms&tbm=isch&sa=X&ved=0ahUKEwjuoPbq2rrQAhWBsxQKHAKsBM4Q\\_AUIBigB#imgsrc=-L14Ymp1v0v04M%3A](https://www.google.pt/search?q=segmentação+por+detecção+de+bordas+e+por+binarização&biw=1366&bih=662&source=lnms&tbm=isch&sa=X&ved=0ahUKEwjuoPbq2rrQAhWBsxQKHAKsBM4Q_AUIBigB#imgsrc=-L14Ymp1v0v04M%3A). [Acessado: 21-Nov-2016].
- [18] O. Raul, A. Valente, O. Prof, M. Eduardo, and R. V. Mesquita, "TÍTULO DO PROJETO : Introdução à Morfologia Matemática Binária e em Tons de Cinza Morfologia Matemática Binária," 2010.
- [19] E. Chevallier, A. Chevallier, and J. Angulo, "N-ary Mathematical Morphology," *Math. Morphol. - Theory Appl.*, vol. 1, no. 1, pp. 42-59, 2016.
- [20] P. Flach, *Machine Learning: the art and science of algorithms that make sense of data*, 1st ed. Cambridge: Cambridge University Press, 2012.
- [21] J.-S. R. Jang, *Neuro-fuzzy and soft computing : a computational approach to learning and machine intelligence*. New Jersey: Upper Saddle River, New Jersey: Prentice

- Hall, 1997.
- [22] L. Breiman, "Randomforest2001," pp. 1-33, 2001.
  - [23] K. J. Wessels *et al.*, "Rapid land cover map updates using change detection and robust random forest classifiers," *Remote Sens.*, vol. 8, no. 11, pp. 1-25, 2016.
  - [24] R. Xinjiang, S. Tian, X. Zhang, J. Tian, and Q. Sun, "Random Forest Classification of Wetland Landcovers from Multi-Sensor Data in the Arid," pp. 1-14, 2016.
  - [25] M. Walker, "Random Forests Algorithm - Data Science Central," 2013. [Online]. Disponível: <http://www.datasciencecentral.com/profiles/blogs/random-forests-algorithm>. [Acessado: 01-Feb-2017].
  - [26] "Gradient boosting," 2017. [Online]. Disponível: [https://en.wikipedia.org/wiki/Gradient\\_boosting](https://en.wikipedia.org/wiki/Gradient_boosting). [Acessado: 01-Feb-2017].
  - [27] R. Blagus and L. Lusa, "Gradient boosting for high-dimensional prediction of rare events," *Comput. Stat. Data Anal.*, 2016.
  - [28] T. Hastie, R. Tibshirani, and J. Friedman, "The Elements of Statistical Learning," *Elements*, vol. 1, pp. 337-387, 2009.
  - [29] G. Napolitano, J. C. Stingl, M. Schmid, and R. Viviani, "Predicting CYP2D6 phenotype from resting brain perfusion images by gradient boosting," *Psychiatry Res. Neuroimaging*, vol. 259, no. July 2016, pp. 16-24, 2017.
  - [30] Y. Bengio, *Learning Deep Architectures for AI*, vol. 2, no. 1. 2009.
  - [31] L. Deng and D. Yu, "Deep Learning: Methods and Applications," *Found. Trends® Signal Process.*, vol. 7, no. 3-4, pp. 197--387, 2013.
  - [32] "Deep learning," 2017. [Online]. Disponível: [https://en.wikipedia.org/wiki/Deep\\_learning](https://en.wikipedia.org/wiki/Deep_learning). [Acessado: 02-Feb-2017].
  - [33] "Machine Learning in Modern Medicine with Erin LeDell at Stanford Med," 2015. [Online]. Disponível: <https://www.slideshare.net/Oxdata/machine-learning-in-modern-medicine-with-erin-ledell-at-stanford-med>. [Acessado: 26-Apr-2017].
  - [34] J. W. Kim, "Classification with Deep Belief Networks."
  - [35] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," *Proc. 26th Annu. Int. Conf. Mach. Learn. ICML 09*, vol. 2008, pp. 1-8, 2009.
  - [36] B. Ribeiro, N. Lopes, and J. Goncalves, "Restricted Boltzmann Machines and Deep Belief Networks on Multi-Core Processors," *Wcci-Ijcnn*, 2012.
  - [37] G. Camps-Valls, D. Tuia, and L. Bruzzone, "Advances in Hyperspectral Image Classification," *IEEE Signal Processing Magazine*, p. 160, 2014.
  - [38] "Understanding Support Vector Machine algorithm from examples (along with code)," 2015. [Online]. Disponível: <https://www.analyticsvidhya.com/blog/2015/10/understaing-support-vector-machine-example-code/>. [Acessado: 26-Apr-2017].
  - [39] C. M. Bishop, *Pattern Recognition and Machine Learning*, vol. 4, no. 4. 2006.
  - [40] K. P. Singh, N. Basant, and S. Gupta, *Support vector machines in water quality management*, vol. 703, no. 2. 2011.
  - [41] J. S. Cardoso and J. F. Pinto, "Learning to Classify Ordinal Data : The Data Replication Method," *J. Mach. Learn. Res.*, vol. 8, pp. 1393-1429, 2007.
  - [42] K. Y. Chang, C. S. Chen, and Y. P. Hung, "Ordinal hyperplanes ranker with cost sensitivities for age estimation," *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern Recognit.*, pp. 585-592, 2011.
  - [43] F. Maselli, A. Rodolfi, L. Bottai, S. Romanelli, and C. Conese, "Classification of Mediterranean vegetation by TM and ancillary data for the evaluation of fire risk," *Int. J. Remote Sens.*, vol. 21, no. 17, pp. 3303-3313, 2000.
  - [44] F. Maselli, S. Romanelli, L. Bottai, and G. Zipoli, "Use of NOAA-AVHRR NDVI images for the estimation of dynamic fire risk in Mediterranean areas," *Remote Sens. Environ.*, vol. 86, no. 2, pp. 187-197, 2003.
  - [45] S. Kuntz and M. Karteris, "Fire Risk Modelling Based on Satellite Remote Sensing and GIS," *EARSeL Advances in Remote Sensing*, vol. 4, no. 3. pp. 39-46, 1995.
  - [46] REDD+ Moçambique, "Florestas em Moçambique." [Online]. Disponível: <http://www.redd.org.mz/page.php?id=55>. [Acessado: 01-Feb-2017].
  - [47] M. P. Temudo and J. M. N. Silva, "Agriculture and forest cover changes in post-war Mozambique," *J. Land Use Sci.*, vol. 4248, no. September, pp. 1-18, 2011.

- [48] T. M. A. Santos, A. Mora, R. A. Ribeiro, and J. M. N. Silva, "Fuzzy-fusion approach for land cover classification," *INES 2016 - 20th Jubil. IEEE Int. Conf. Intell. Eng. Syst. Proc.*, pp. 177-182, 2016.
- [49] "Statistics - Tuning Parameter [Gerardnico]," 2017. [Online]. Disponível: [https://gerardnico.com/wiki/data\\_mining/tuning\\_parameter](https://gerardnico.com/wiki/data_mining/tuning_parameter). [Acessado: 04-May-2017].
- [50] T. G. Dietterich, "Ensemble Methods in Machine Learning," *Mult. Classif. Syst.*, vol. 1857, pp. 1-15, 2000.
- [51] "Planet: Understanding the Amazon from Space | Kaggle," 2017. [Online]. Disponível: <https://www.kaggle.com/c/planet-understanding-the-amazon-from-space>. [Acessado: 16-May-2017].
- [52] B. Krawczyk, L. L. Minku, J. Gama, J. Stefanowski, and M. Woźniak, "Ensemble learning for data stream analysis: A survey," *Inf. Fusion*, vol. 37, pp. 132-156, Sep. 2017.
- [53] I. Žliobaitė, M. Pechenizkiy, and J. Gama, "An Overview of Concept Drift Applications," *Big Data Anal. New Algorithms a New Soc.*, vol. 16, pp. 91-114, 2016.
- [54] "How To Compare Machine Learning Algorithms in Python with scikit-learn - Machine Learning Mastery," 2016. [Online]. Disponível: <http://machinelearningmastery.com/compare-machine-learning-algorithms-python-scikit-learn/>. [Acessado: 16-May-2017].
- [55] MINISTÉRIO PARA A COORDENAÇÃO DA ACÇÃO AMBIENTAL, "Avaliação da vulnerabilidade as mudanças climáticas e estratégias de adaptação," p. 61, 2005.